

# Hierarchical models

Dr. Jarad Niemi

Iowa State University

August 29, 2017

# Normal hierarchical model

Let

$$Y_{ij} \stackrel{\text{ind}}{\sim} N(\theta_i, \sigma^2)$$

for  $i = 1, \dots, I$ ,  $j = 1, \dots, n_i$ , and  $\sum_{i=1}^I n_i = n$ . Now consider the following model assumptions:

- $\theta_i \stackrel{\text{ind}}{\sim} N(\mu, \tau^2)$
- $\theta_i \stackrel{\text{ind}}{\sim} La(\mu, \tau)$
- $\theta_i \stackrel{\text{ind}}{\sim} t_\nu(\mu, \tau^2)$
- $\theta_i \stackrel{\text{ind}}{\sim} \pi\delta_0 + (1 - \pi)N(\mu, \tau^2)$
- $\theta_i \stackrel{\text{ind}}{\sim} \pi\delta_0 + (1 - \pi)t_\nu(\mu, \tau^2)$

To perform a Bayesian analysis, we need a prior on  $\mu$ ,  $\tau^2$ , and (in the case of the discrete mixture)  $\pi$ .

# Normal hierarchical model

Consider the model

$$\begin{aligned} Y_{ig} &\stackrel{\text{ind}}{\sim} N(\theta_g, \sigma^2) \\ \theta_g &\stackrel{\text{ind}}{\sim} N(\mu, \tau^2) \end{aligned}$$

where  $i = 1, \dots, n_g$ ,  $g = 1, \dots, G$ , and  $n = \sum_{g=1}^G n_g$  with prior distribution

$$p(\mu, \sigma^2, \tau^2) = p(\mu)p(\sigma^2)p(\tau^2) \propto \frac{1}{\sigma^2} Ca^+(\tau; 0, C).$$

For background on why we are using these priors for the variances, see Gelman (2006) <https://projecteuclid.org/euclid.ba/1340371048>: “Prior distributions for variance parameters in hierarchical models (comment on article by Browne and Draper)”.

# Gibbs sampler for normal hierarchical model

Here is a possible Gibbs sampler for this model:

- For  $g = 1, \dots, G$ , sample  $\theta_g \sim p(\theta_g | \dots)$ .
- Sample  $\sigma^2 \sim p(\sigma^2 | \dots)$ .
- Sample  $\mu \sim p(\mu | \dots)$ .
- Sample  $\tau^2 \sim p(\tau^2 | \dots)$ .

How many steps exist in this Gibbs sampler?  $G+3$ ? 4?

## 2-Step Gibbs sampler for normal hierarchical model

Here is a 2-step Gibbs sampler:

1. Sample  $\theta = (\theta_1, \dots, G) \sim p(\theta | \dots)$ .
2. Sample  $\mu, \sigma^2, \tau^2 \sim p(\mu, \sigma^2, \tau^2 | \dots)$ .

There is stronger theoretical support for 2-step Gibbs sampler, thus, if we can, it is prudent to construct a 2-step Gibbs sampler.

# Sampling $\theta$

The full conditional for  $\theta$  is

$$\begin{aligned} p(\theta | \dots) &\propto p(\theta, \mu, \sigma^2, \tau^2 | y) \\ &\propto p(y | \theta, \sigma^2) p(\theta | \mu, \tau^2) p(\mu, \sigma^2, \tau^2) \\ &\propto p(y | \theta, \sigma^2) p(\theta | \mu, \tau^2) \\ &= \prod_{g=1}^G p(y_g | \theta_g, \sigma^2) p(\theta_g | \mu, \tau^2) \end{aligned}$$

where  $y_g = (y_{1g}, \dots, y_{n_g g})$ . We now know that the  $\theta_g$  are conditionally independent of each other.

# Sampling $\theta_g$

The full conditional for  $\theta_g$  is

$$\begin{aligned} p(\theta_g | \cdots) &\propto p(y_g | \theta_g, \sigma^2) p(\theta_g | \mu, \tau^2) \\ &= \prod_{i=1}^{n_g} N(y_{ig}; \theta_g, \sigma^2) N(\theta_g; \mu, \tau^2) \end{aligned}$$

Notice that this does not include  $\theta_{g'}$  for any  $g' \neq g$ . This is an alternative way to conclude that the  $\theta_g$  are conditionally independent of each other.

Thus

$$\theta_g | \cdots \stackrel{ind}{\sim} N(\mu_g, \tau_g^2)$$

where

$$\begin{aligned} \tau_g^2 &= [\tau^{-2} + n_g \sigma^{-2}]^{-1} \\ \mu_g &= \tau_g^2 [\mu \tau^{-2} + \bar{y}_g n_g \sigma^{-2}] \\ \bar{y}_g &= \frac{1}{n_g} \sum_{i=1}^{n_g} y_{ig}. \end{aligned}$$

# Sampling $\mu, \sigma^2, \tau^2$

The full conditional for  $\mu, \sigma^2, \tau^2$  is

$$\begin{aligned} p(\mu, \sigma^2, \tau^2 | \dots) &\propto p(y|\theta, \sigma^2)p(\theta|\mu, \tau^2)p(\mu)p(\sigma^2)p(\tau^2) \\ &= p(y|\theta, \sigma^2)p(\sigma^2)p(\theta|\mu, \tau^2)p(\mu)p(\tau^2) \end{aligned}$$

So we know that  $\sigma^2$  is independent of  $\mu$  and  $\tau^2$ .



# Sampling $\sigma^2$

Recall that

$$y_{ig} \stackrel{\text{ind}}{\sim} N(\theta_g, \sigma^2) \text{ and } p(\sigma^2) \propto 1/\sigma^2.$$

Thus, we are in the scenario of normal data with a known mean and unknown variance and the unknown variance has our default prior. Thus, we should know the full conditional is

$$\sigma^2 | \dots \sim IG\left(\frac{n}{2}, \frac{1}{2} \sum_{g=1}^G \sum_{i=1}^{n_g} (y_{ig} - \theta_g)^2\right).$$

To derive the full conditional, use

$$\begin{aligned} p(\sigma^2 | \dots) &\propto \prod_{g=1}^G \prod_{i=1}^{n_g} (\sigma^2)^{-1/2} \exp\left(-\frac{1}{2\sigma^2} (y_{ig} - \theta_g)^2\right) \frac{1}{\sigma^2} \\ &= (\sigma^2)^{-n/2-1} \exp\left(-\frac{1}{2} \sum_{g=1}^G \sum_{i=1}^{n_g} (y_{ig} - \theta_g)^2 / \sigma^2\right) \end{aligned}$$

which is the kernel of a  $IG\left(\frac{n}{2}, \frac{1}{2} \sum_{g=1}^G \sum_{i=1}^{n_g} (y_{ig} - \theta_g)^2\right)$ .

# Sampling $\mu, \tau^2$

Recall that

$$\theta_g \stackrel{\text{ind}}{\sim} N(\mu, \tau^2) \text{ and } p(\mu, \tau^2) \propto Ca^+(\tau; 0, C).$$

This is a non-standard distribution, but is extremely close a normal model with unknown mean and variance with the standard non-informative prior  $p(\mu, \tau^2) \propto 1/\tau^2$  or the conjugate normal-inverse-gamma prior.

Here are some options for sampling from this distribution:

- random-walk Metropolis (in 2 dimensions),
- independent Metropolis-Hastings using posterior from standard non-informative prior as the proposal, or
- rejection sampling using posterior from standard non-informative prior as the proposal

The posterior under the standard non-informative prior is

$$\tau^2 | \dots \sim \text{Inv-}\chi^2(G-1, s_\theta^2) \text{ and } \mu | \tau^2, \dots \sim N(\bar{\theta}, \tau^2/G)$$

where  $\bar{\theta} = \frac{1}{G} \sum_{g=1}^G \theta_g$  and  $s_\theta^2 = \frac{1}{G-1} (\sum_{g=1}^G \theta_g^2 - G\bar{\theta}^2)$ . What is the MH ratio?

# Markov chain Monte Carlo for normal hierarchical model

1. Sample  $\theta \sim p(\theta | \dots)$ :
  - a. For  $g = 1, \dots, G$ , sample  $\theta_g \sim N(\mu_g, \tau_g^2)$ .
2. Sample  $\mu, \sigma^2, \tau^2$ :
  - a. Sample  $\sigma^2 \sim IG(n/2, SSE)$ .
  - b. Sample  $\mu, \tau^2$  using independent Metropolis-Hastings using posterior from standard non-informative prior as the proposal.

What happens if  $\theta_g \stackrel{ind}{\sim} La(\mu, \tau)$  or  $\theta_g \stackrel{ind}{\sim} t_\nu(\mu, \tau^2)$ ?

# Scale mixtures of normals

Recall that if

$$\theta|\phi \sim N(\phi, V) \text{ and } \phi \sim N(m, C)$$

then

$$\theta \sim N(m, V + C).$$

This is called a location mixture.

Now, if

$$\theta|\phi \sim N(m, C\phi)$$

and we assume a mixing distribution for  $\phi$ , we have a scale mixture. Since the top level distributional assumption is normal, we refer to this as a **scale mixture of normals**.

# t distribution

Let

$$\theta|\phi \sim N(m, \phi C) \text{ and } \phi \sim IG(a, b)$$

then

$$\begin{aligned} p(\theta) &= \int p(\theta|\phi)p(\phi)d\phi \\ &= (2\pi\sqrt{C})^{-1/2} \frac{b^a}{\Gamma(a)} \int \phi^{-1/2} e^{-(\theta-m)^2/2\phi C} \phi^{-(a+1)} e^{-b/\phi} d\phi \\ &= (2\pi C)^{-1/2} \frac{b^a}{\Gamma(a)} \int \phi^{-(a+1/2+1)} e^{-[b+(\theta-m)^2/2C]/\phi} d\phi \\ &= (2\pi C)^{-1/2} \frac{b^a}{\Gamma(a)} \frac{\Gamma(a+1/2)}{[b+(\theta-m)^2/2C]^{a+1/2}} \\ &= \frac{\Gamma([2a+1]/2)}{\Gamma(2a/2)\sqrt{2a\pi bC/a}} \left[ 1 + \frac{1}{2a} \frac{(\theta-m)^2}{bC/a} \right]^{-[2a+1]/2} \end{aligned}$$

Thus

$$\theta \sim t_{2a}(m, bC/a)$$

i.e.  $\theta$  has a  $t$  distribution with  $2a$  degrees of freedom, location  $m$ , scale  $bC/a$ , and variance  $\frac{bC}{a-1}$ .

# Hierarchical $t$ distribution

Let  $m = \mu$ ,  $C = 1$ ,  $a = \nu/2$ , and  $b = \nu\tau^2/2$ , i.e.

$$\theta|\phi \sim N(\mu, \phi) \text{ and } \phi \sim IG(\nu/2, \nu\tau^2/2).$$

Then, we have

$$\theta \sim t_\nu(\mu, \tau^2),$$

i.e. a  $t$  distribution with  $\nu$  degrees of freedom, location  $\mu$ , and scale  $\tau^2$ .

Notice that the parameterization has a redundancy between  $C$  and  $a/b$ , i.e. we could have chosen  $C = \tau^2$ ,  $a = \nu/2$ , and  $b = \nu/2$  and we would have obtained the same marginal distribution for  $\theta$ .

# Laplace distribution

Let

$$\theta|\phi \sim N(m, \phi C^2) \text{ and } \phi \sim \text{Exp}(1/2b^2)$$

where  $E[\phi] = 2b^2$  and  $\text{Var}[\phi] = 4b^4$ . then, by an extension of equation (4) in Park and Casella (2008), we have

$$p(\theta) = \frac{1}{2Cb} e^{-\frac{|\theta-m|}{Cb}}.$$

This is the pdf for a Laplace (double exponential) distribution with location  $m$  and scale  $Cb$  which we write

$$\theta \sim \text{La}(m, Cb).$$

and say  $\theta$  has a Laplace distribution with location  $m$  and scale  $Cb$  and  $E[\theta] = m$  and  $\text{Var}[\theta] = 2[Cb]^2 = 2C^2b^2$ .

# Hierarchical Laplace distribution

Let  $m = \mu$ ,  $C = 1$ , and  $b = \tau$  i.e.

$$\theta|\phi \sim N(\mu, \phi) \text{ and } \phi \sim \text{Exp}(1/2\tau^2).$$

Then, we have

$$\theta \sim \text{La}(\mu, \tau),$$

i.e. a Laplace distribution with location  $\mu$  and scale  $\tau$ .

Notice that the parameterization has a redundancy between  $C$  and  $b$ , i.e. we could have chosen  $C = \tau$  and  $b = 1$  and we would have obtained the same marginal distribution for  $\theta$ .



# Normal hierarchical model

Recall our hierarchical model

$$Y_{ij} \stackrel{ind}{\sim} N(\theta_i, \sigma^2)$$

for  $i = 1, \dots, I$  and  $j = 1, \dots, n_i$ . Now consider the following model assumptions:

- $\theta_i \stackrel{ind}{\sim} N(\mu, \phi_i), \phi_i = \tau^2 \implies \theta_i \stackrel{ind}{\sim} N(\mu, \tau^2)$
- $\theta_i | \phi_i \stackrel{ind}{\sim} N(\mu, \phi_i), \phi_i \stackrel{ind}{\sim} \text{Exp}(1/2\tau^2) \implies \theta_i \stackrel{ind}{\sim} \text{La}(\mu, \tau)$
- $\theta_i | \phi_i \stackrel{ind}{\sim} N(\mu, \phi_i), \phi_i \stackrel{ind}{\sim} \text{IG}(v/2, v\tau^2/2) \implies \theta_i \stackrel{ind}{\sim} t_v(\mu, \tau^2)$

For simplicity, let's assume  $\sigma^2 \sim \text{IG}(a, b)$ ,  $\mu \sim N(m, C)$ , and  $\tau \sim \text{Ca}^+(0, c)$  and that  $\sigma^2$ ,  $\mu$ , and  $\tau$  are *a priori* independent.

# Gibbs sampling

The following Gibbs sampler will converge to the posterior  $p(\theta, \sigma, \mu, \phi, \tau|y)$ :

1. Sample  $\mu \sim p(\mu|\cdots)$ .
2. Independently, sample  $\theta_i \sim p(\theta_i|\cdots)$ .
3. Sample  $\sigma \sim p(\sigma|\cdots)$ .
4. Independently, sample  $\phi_i \sim p(\phi_i|\cdots)$ .
5. Sample  $\tau \sim p(\tau|\cdots)$ .

The first three steps will be common to all models while the last two steps will be unique to each model (without a point mass).

# Sample $\mu$

$$\theta_i \stackrel{ind}{\sim} N(\mu, \phi_i) \text{ and } \mu \sim N(m, C)$$

Immediately, we should know that

$$\mu | \cdots \sim N(m', C')$$

with

$$\begin{aligned} C' &= \left( \frac{1}{C} + \sum_{i=1}^I \frac{1}{\phi_i} \right)^{-1} \\ m' &= C' \left( \frac{m}{C} + \sum_{i=1}^I \frac{\theta_i}{\phi_i} \right) \end{aligned}$$

# Sample $\theta$

$$Y_{ij} \stackrel{\text{ind}}{\sim} N(\theta_i, \sigma^2) \text{ and } \theta_i \sim N(\mu, \phi_i)$$

$$\begin{aligned} p(\theta | \dots) &\propto \left[ \prod_{i=1}^I \prod_{j=1}^{n_i} e^{-(y_{ij} - \theta_i)^2 / 2\sigma^2} \right] \left[ \prod_{i=1}^I e^{-(\theta_i - \mu)^2 / 2\phi_i} \right] \\ &\propto \prod_{i=1}^I \left[ \prod_{j=1}^{n_i} e^{-(y_{ij} - \theta_i)^2 / 2\sigma^2} e^{-(\theta_i - \mu)^2 / 2\phi_i} \right] \end{aligned}$$

Thus  $\theta_i$  are conditionally independent given everything else. It should be obvious that

$$\theta_i | \dots \sim N \left( \left[ \frac{\mu}{\phi_i} + \frac{n_i}{\sigma^2} \bar{y}_i \right], \left[ \frac{1}{\phi_i} + \frac{n_i}{\sigma^2} \right]^{-1} \right)$$

where  $\bar{y}_i = \sum_{j=1}^{n_i} y_{ij} / n_i$ .

## Sample $\sigma^2$

$$Y_{ij} \stackrel{\text{ind}}{\sim} N(\theta_i, \sigma^2) \text{ and } \sigma^2 \sim IG(a, b)$$

This is just a normal data model with an unknown variance that has the conjugate prior. The only difficulty is that we have several groups here. But very quickly you should be able to determine that

$$\sigma^2 | \dots \sim IG(a', b')$$

where

$$\begin{aligned} a' &= a + \sum_{i=1}^I n_i / 2 = a + n / 2 \\ b' &= b + \sum_{i=1}^I \sum_{j=1}^{n_i} (y_{ij} - \theta_i)^2 / 2. \end{aligned}$$

## Distributional assumption for $\theta_i$

$$Y_{ij} \stackrel{\text{ind}}{\sim} N(\theta_i, \sigma^2) \text{ and } \theta_i \stackrel{\text{ind}}{\sim} N(\mu, \phi_i)$$

$$\phi_i = \tau$$

$$\phi_i \sim \text{Exp}(1/2\tau^2)$$

$$\phi_i \sim \text{IG}(v/2, v\tau^2/2)$$

The steps that are left are 1) sample  $\phi$  and 2) sample  $\tau^2$ ,

# Sample $\phi$ for normal model

For normal model,  $\phi_i = \tau$ , so we will address this when we sample  $\tau$ .

## Sample $\phi$ for Laplace model

For Laplace model,

$$\theta_i \stackrel{\text{ind}}{\sim} N(\mu, \phi_i) \text{ and } \phi_i \stackrel{\text{ind}}{\sim} \text{Exp}(1/2\tau^2),$$

so the full conditional is

$$p(\phi | \dots) \propto \left[ \prod_{i=1}^I N(\theta_i; \mu, \phi_i) \text{Exp}(\phi_i; 1/2\tau^2) \right].$$

So the individual  $\phi_i$  are conditionally independent with

$$p(\phi_i | \dots) \propto N(\theta_i; \mu, \phi_i) \text{Exp}(\phi_i; 1/2\tau^2) \propto \phi_i^{-1/2} e^{-(\theta_i - \mu)^2 / 2\phi_i} e^{-\phi_i / 2\tau^2}$$

If we perform the transformation  $\eta_i = 1/\phi_i$ , we have

$$p(\eta_i | \dots) \propto \eta_i^{-3/2} e^{-\frac{(\theta_i - \mu)^2}{2} \eta_i - \frac{1}{2\tau^2 \eta_i}}$$

which is the kernel of an inverse Gaussian distribution with mean  $\sqrt{1/\tau^2(\theta_i - \mu)^2}$  and scale  $1/\tau^2$  where the parameterization is such that the variance is  $\mu^3/\lambda$  (different from the `mgcv::rig` parameterization).



## Sample $\phi$ for $t$ model

For the  $t$  model,

$$\theta_i \stackrel{\text{ind}}{\sim} N(\mu, \phi_i) \text{ and } \phi_i \stackrel{\text{ind}}{\sim} IG(v/2, v\tau^2/2),$$

so we have

$$\phi_i | \dots \stackrel{\text{ind}}{\sim} IG([v+1]/2, [v\tau^2 + (\theta_i - \mu)^2]/2).$$

Since this is just  $I$  independent normal data models with a known mean and independent conjugate inverse gamma priors on the variance.

## Sample $\tau$ for normal model

Let

$$\theta_i \stackrel{\text{ind}}{\sim} N(\mu, \tau^2) \text{ and } \tau \sim Ca^+(0, c).$$

so the full conditional is

$$p(\eta | \dots) \propto \eta^{-I/2} e^{-\sum_{i=1}^I (\theta_i - \mu)^2 / 2\eta} (1 + \eta/c^2)^{-1} \eta^{-1/2}$$

where we performed the transformation  $\eta = \tau^2$  on the prior.

Let's use Metropolis-Hastings with proposal distribution

$$IG\left(\frac{I-1}{2}, \sum_{i=1}^I \frac{(\theta_i - \mu)^2}{2}\right)$$

and acceptance probability  $\min\{1, \rho\}$  where

$$\rho = \frac{(1 + \eta^*/c^2)^{-1}}{(1 + \eta^{(i)}/c^2)^{-1}} = \frac{1 + \eta^{(i)}/c^2}{1 + \eta^*/c^2}$$

where  $\eta^{(i)}$  and  $\eta^*$  are the current and proposed value respective.

## Sample $\tau$ for Laplace model

Let

$$\phi_i \sim \text{Exp}(1/2\tau^2) \text{ and } \tau \sim \text{Ca}^+(0, c)$$

so the full conditional is

$$p(\eta | \dots) \propto \eta^{-I} e^{-\sum_{i=1}^I \phi_i / 2\eta} (1 + \eta/c^2)^{-1} \eta^{-1/2}.$$

Let's use Metropolis-Hastings with proposal distribution

$$IG\left(I - \frac{1}{2}, \sum_{i=1}^I \frac{\phi_i}{2}\right)$$

and acceptance probability  $\min\{1, \rho\}$  where again

$$\rho = \frac{1 + \eta^{(i)}/c^2}{1 + \eta^*/c^2}.$$

Then we calculate  $\tau = \sqrt{\eta}$ .

## Sample $\tau$ for $t$ model

Let

$$\phi_i \sim IG(v/2, v\tau^2/2) \text{ and } \tau \sim Ca^+(0, c)$$

so the full conditional is

$$p(\eta | \dots) \propto \eta^{Iv/2} e^{-\frac{\eta}{2} \sum_{i=1}^I \frac{1}{\phi_i}} (1 + \eta/c^2)^{-1} \eta^{-1/2}.$$

Let's use Metropolis-Hastings with proposal distribution

$$Ga\left(\frac{Iv+1}{2}, \frac{1}{2} \sum_{i=1}^I \frac{1}{\phi_i}\right)$$

and acceptance probability  $\min\{1, \rho\}$  where again

$$\rho = \frac{1 + \eta^{(i)}/c^2}{1 + \eta^*/c^2}.$$

Then we calculate  $\tau = \sqrt{\eta}$ .

## Dealing with point-mass distributions

We would also like to consider models with

$$\theta_i \stackrel{\text{ind}}{\sim} \pi \delta_0 + (1 - \pi) N(\mu, \phi_i)$$

where  $\phi_i = \tau^2$  corresponds to a normal and

$$\phi_i \stackrel{\text{ind}}{\sim} IG(v/2, v\tau^2/2)$$

corresponds to a  $t$  distribution for the non-zero  $\theta_i$ .

Similar to the previous, the  $\theta_i$  are conditionally independent. To sample  $\theta_i$ , we calculate

$$\begin{aligned} \pi' &= \frac{\pi \prod_{j=1}^{n_i} N(y_{ij}; 0, \sigma^2)}{\pi \prod_{j=1}^{n_i} N(y_{ij}; 0, \sigma^2) + (1 - \pi) \prod_{j=1}^{n_i} N(y_{ij}; \mu, \phi_i + \sigma^2)} \\ \phi'_i &= \left( \frac{1}{\phi_i} + \frac{n_i}{\sigma^2} \right)^{-1} \\ \mu'_i &= \phi'_i \left( \frac{\mu}{\phi_i} + \frac{n_i}{\sigma^2} \bar{y}_i \right) \end{aligned}$$

# Dealing with point-mass distributions (cont.)

Let

$$\theta_i \overset{\text{ind}}{\sim} \pi\delta_0 + (1 - \pi)N(\mu, \phi_i)$$

and independently  $\pi \sim \text{Beta}(s, f)$ ,  $\mu \sim N(m, C)$ , and  $\phi_i = \tau^2$  for normal model or  $\phi_i \overset{\text{ind}}{\sim} \text{IG}(v/2, v\tau^2/2)$  for the  $t$  model.

The full conditional for  $\pi$  is

$$\pi | \dots \sim \text{Beta} \left( s + \sum_{i=1}^I \text{I}(\theta_i = 0), f + \sum_{i=1}^I \text{I}(\theta_i \neq 0) \right)$$

and  $\mu$  and  $\phi_i$  get updated using only those  $\theta_i$  that are non-zero.