

Class 5: Multivariate and multi-layer hierarchical models

Andrew Parnell
andrew.parnell@ucd.ie



Learning outcomes:

- ▶ Understand how to add in multiple layers to a hierarchical model
- ▶ Follow a detailed example of building a model
- ▶ Be able to fit multivariate models in JAGS
- ▶ Be able to work with missing data in JAGS

1 / 40

2 / 40

Some new terminology

- ▶ Most of the models we have covered so far contain only one hidden or *latent* set of parameters
- ▶ For example, the data y may depend on a parameter β , which itself depends on a parameter θ . θ is given a prior distribution
- ▶ We say that the data are at the 'top level', the parameter β is a *latent parameter* at the second level, and the hyper-parameter θ is also a latent parameter at the third level
- ▶ We say that the prior distribution on β is *conditional* on θ , whilst the prior distribution (if it just involves numbers) is a *marginal prior distribution*

What is a multi-layer model?

- ▶ A multi-layer model is one where we have many (usually more than 2 or 3) layers of parameters conditional on each other
- ▶ It's very straightforward to add in these extra layers in JAGS/Stan
- ▶ The question is whether they are necessary or not, and how much the data can tell us about them

3 / 40

4 / 40

A simple example

- ▶ We will work through the earnings example, extending it to produce a much more complex model
- ▶ The model we last met had the following JAGS code:

```
,  
...  
y[i] ~ dnorm(alpha[eth[i]] +  
             beta[eth[i]]*(x[i] - mean(x)),  
             sigma^-2)  
}  
...  
# Priors  
for(j in 1:N_eth) {  
  alpha[j] ~ dnorm(mu_alpha, sigma_alpha^-2)  
  beta[j] ~ dnorm(mu_beta, sigma_beta^-2)  
},
```

5 / 40

Adding extra layers

- ▶ The priors for the hyper-parameters at the bottom of the model were

```
,  
...  
mu_alpha ~ dnorm(11, 2^-2)  
mu_beta ~ dnorm(0, 0.1^-2)  
sigma ~ dunif(0, 5)  
sigma_alpha ~ dunif(0, 2)  
sigma_beta ~ dunif(0, 2)  
},
```

- ▶ There is nothing to stop us adding in extra parameters beneath these layers

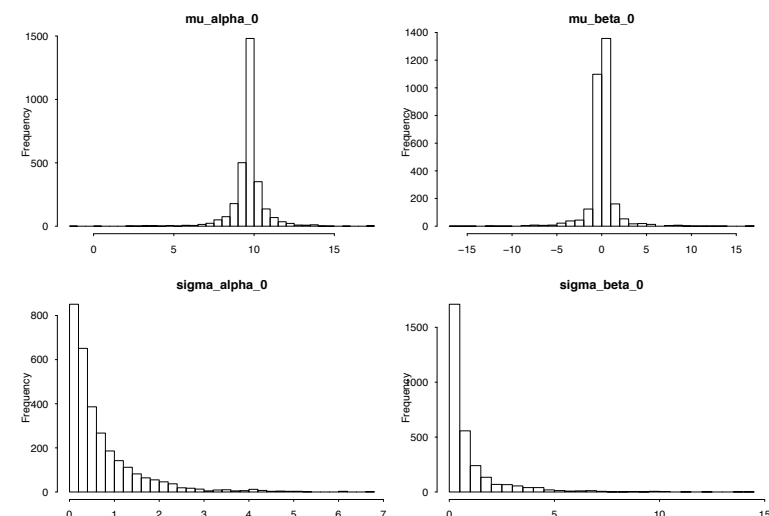
6 / 40

New JAGS code

```
jags_code = '  
model{  
  # Likelihood  
  for(i in 1:N) {  
    y[i] ~ dnorm(alpha[eth[i]] +  
                beta[eth[i]]*(x[i] - mean(x)),  
                sigma^-2)  
  }  
  # Priors  
  for(j in 1:N_eth) {  
    alpha[j] ~ dnorm(mu_alpha, sigma_alpha^-2)  
    beta[j] ~ dnorm(mu_beta, sigma_beta^-2)  
  }  
  mu_alpha ~ dnorm(mu_alpha_0, sigma_alpha_0^-2)  
  mu_alpha_0 ~ dnorm(0, 20^-2)  
  sigma_alpha_0 ~ dt(0,5,1)T(0,)  
  mu_beta ~ dnorm(mu_beta_0, sigma_beta_0^-2)  
  mu_beta_0 ~ dnorm(0, 20^-2)  
  sigma_beta_0 ~ dt(0,5,1)T(0,)  
  sigma ~ dunif(0, 5)  
  sigma_alpha ~ dt(0,5,1)T(0,)  
  sigma_beta ~ dt(0,5,1)T(0,)  
}  
,
```

7 / 40

What information is there about these extra parameters?



8 / 40

How many layers should I add?

- ▶ We could go on adding layers here if we wanted to, but it's not clear what benefit it would have
- ▶ For example, do we really need to know the standard deviation of the mean of the intercept? The answer to this will depend on the questions of interest, and the amount and type of prior information available
- ▶ A general rule is, in the absence of any strong prior information, to add one extra layer of parameters beyond the latent parameters of key interest

9 / 40

Writing the same model in different ways

- ▶ The model on the previous slides has a different intercept and slope for each ethnicity group, with the information about them tied together through the prior distributions on them
- ▶ The likelihood was written as:

```
y[i] ~ dnorm(alpha[eth[i]] +  
             beta[eth[i]]*(x[i] - mean(x)),  
             sigma^-2)
```

which in maths can be written as:

$$y_i \sim N(\alpha_{\text{eth}_i} + \beta_{\text{eth}_i} x_i, \sigma^2)$$

where eth_i takes the values 1, 2, 3, or 4

- ▶ Remember, y_i is the log-earnings of individual i where $i = 1, \dots, N$

10 / 40

Re-writing the model

- ▶ Commonly you'll see y here re-defined as y_{ij} where $j = 1, \dots, 4$ represents ethnicity, and $i = 1, \dots, N_j$ is the number of individuals with ethnicity j
- ▶ The likelihood can then be written as:

$$y_{ij} \sim N(\alpha_j + \beta_j x_i, \sigma^2)$$

- ▶ Note that this is exactly the same model, just re-written slightly differently. In fact, this latter model is much harder to write out in JAGS/Stan code

11 / 40

Fixed vs random effect models

- ▶ Thinking about this model in more detail

$$y_{ij} \sim N(\alpha_j + \beta_j x_i, \sigma^2)$$

- ▶ If the α_j and β_j parameters are all given independent prior distributions, e.g. $\alpha_j \sim N(0, 100)$ or similar, then this is considered a *fixed effects* model
- ▶ If the α_j and β_j are given prior distributions that tie the values together, e.g. $\alpha_j \sim N(\mu_\alpha, \sigma_\alpha^2)$, then this is often called a *random effects* model
- ▶ (In fact, nobody can agree on what a fixed or random effects model actually is)
- ▶ If you have multiple different types of priors then this is called a *mixed effects* model

12 / 40

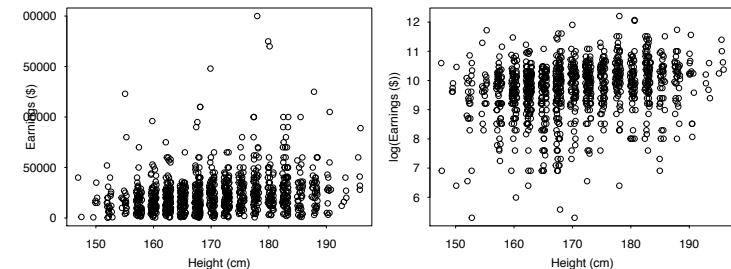
Mixed effects vs hierarchical models

- ▶ The hierarchical models we have been studying all use the *random effects* approach wherever possible
- ▶ The big advantage of using this approach is that we get to *borrow* strength between the different groups (here eth, but it could be anything)
- ▶ Whenever we have a categorical covariate we should always be putting a constraining/tying prior distribution on them, and looking at how the effects vary between the groups
- ▶ This is strongly linked to the idea of *partial pooling* which we will meet later in the course

13 / 40

Example: multi-layer earnings data

- ▶ We will now go through and build a much more complicated model for the earnings data, taken from the Gelman and Hill book, using only weak priors
- ▶ We can generate data from these models (either using the prior or the posterior), and we can draw a DAG
- ▶ Our goal is to explore the factors which explain earnings. We have variables on height, age, and ethnicity.
- ▶ If we first plot the data



14 / 40

Transformations

- ▶ From the left-hand plot there seem to be quite a few extreme observations, and there's a possibility that the relationship between height and earnings is non-linear
- ▶ The right-hand plot seems to have stabilised most of the extreme observations, and perhaps linearity is more appropriate
- ▶ Notice that a linear model implies:

$$y_i = \alpha + \beta x_i + \epsilon_i$$

whilst the log-linear model implies:

$$y_i = \exp(\alpha + \beta x_i + \epsilon_i) = e^\alpha \times e^{\beta x_i} \times e^{\epsilon_i}$$

so the coefficients, once exponentiated, have multiplicative effects that are relatively easy to interpret

15 / 40

Fitting the first model

- ▶ If we fit a model with just height (mean centered) we get the following JAGS output

```
## Inference for Bugs model at "6", fit using jags,
## 3 chains, each with 2000 iterations (first 1000 discarded)
## n.sims = 3000 iterations saved
##               mu.vect sd.vect   2.5%   25%   50%
## alpha         9.738   0.028   9.684   9.719   9.738
## beta_height    0.023   0.003   0.017   0.021   0.023
## sigma         0.906   0.019   0.869   0.893   0.906
## deviance      2798.491   2.388  2795.748  2796.737  2797.819
##               Rhat n.eff
## alpha         1.002  1200
## beta_height    1.001  3000
## sigma         1.001  3000
## deviance      1.001  3000
##
## For each parameter, n.eff is a crude measure of effective
## and Rhat is the potential scale reduction factor (at convergence)
```

16 / 40

Interpreting the parameters

- ▶ These parameters are directly interpretable:
 - ▶ The mean of the log earnings at the mean height is about 9.737, which is about 17k on the original scale
 - ▶ We can also get e.g. a 95% confidence interval using the JAGS output. From 16000 to 18000
 - ▶ For every extra cm so you gain 0.0225 on the log scale, i.e. an 2.28% gain in income
 - ▶ From the posterior of σ , we can guess that about 68% of predictions will be within 0.906 on the log scale or within a factor of about 2.48 of the prediction
- ▶ Interpretation for the intercept would have been harder had we not mean-centered the height variable
- ▶ The DIC is 2801.34 with 2.85 effective parameters

17 / 40

Improving the model

- Now suppose we fit a model with a random intercept for ethnicity

```
## Inference for Bugs model at "7", fit using jags,
## 3 chains, each with 2000 iterations (first 1000 discarded)
## n.sims = 3000 iterations saved
##          mu.vect sd.vect   2.5%    25%    50%    75%   97.5%
## alpha      9.710   0.087   9.520   9.669   9.716   9.755   9.872
## beta_eth[1] -0.022   0.090  -0.232  -0.061  -0.011   0.015   0.161
## beta_eth[2] -0.043   0.100  -0.284  -0.082  -0.023   0.008   0.123
## beta_eth[3]  0.039   0.087  -0.126  -0.003   0.025   0.080   0.232
## beta_eth[4]  0.030   0.110  -0.164  -0.022   0.010   0.074   0.295
## beta_height 0.022   0.003   0.017   0.020   0.022   0.024   0.028
## sigma      0.907   0.020   0.868   0.893   0.906   0.920   0.947
## deviance   2797.911   2.991 2793.774 2795.766 2797.367 2799.402 2805.235
##          Rhat n.eff
## alpha      1.002  3000
## beta_eth[1] 1.001  3000
## beta_eth[2] 1.003  1500
## beta_eth[3] 1.001  3000
## beta_eth[4] 1.004  1500
## beta_height 1.001  3000
## sigma      1.001  3000
## deviance   1.001  3000
##
## For each parameter, n.eff is a crude measure of effective sample size,
## and Rhat is the potential scale reduction factor (at convergence, Rhat=1).
##
## DIC info (using the rule, pD = var(deviance)/2)
## pD = 4.5 and DIC = 2802.4
## DIC is an estimate of expected predictive error (lower deviance is better).
```

18 / 40

Interpreting the output

- ▶ The parameters α and β_{height} haven't changed much in the mean
- ▶ The 95% confidence interval for α has increased: 14000 to 19000
- ▶ The DIC is 2802.39 with 4.47 effective parameters. Pretty much the same as above
- ▶ We also have estimates for each ethnicity, none of these have a strong effect away from zero

19 / 40

Now an interaction model

```
## Inference for Bugs model at "8", fit using jags,
## 3 chains, each with 2000 iterations (first 1000 discarded)
## n.sims = 3000 iterations saved
##          mu.vect sd.vect   2.5%    25%    50%    75%
## alpha      9.699   0.087   9.522   9.658   9.706   9.745
## beta_eth[1] -0.021   0.092  -0.215  -0.059  -0.012   0.019
## beta_eth[2] -0.046   0.102  -0.283  -0.091  -0.027   0.006
## beta_eth[3]  0.046   0.088  -0.109  -0.001   0.034   0.088
## beta_eth[4]  0.019   0.107  -0.180  -0.029   0.007   0.059
## beta_height[1] 0.009   0.009  -0.008   0.003   0.009   0.015
## beta_height[2] 0.013   0.010  -0.008   0.006   0.013   0.019
## beta_height[3] 0.025   0.003   0.019   0.023   0.025   0.027
## beta_height[4] 0.007   0.017  -0.030  -0.002   0.009   0.019
## mu_beta_height 0.013   0.019  -0.026   0.007   0.015   0.021
## sigma      0.904   0.020   0.864   0.891   0.905   0.918
## sigma_height 0.025   0.030   0.003   0.010   0.016   0.027
## deviance   2793.128   3.700 2787.518 2790.414 2792.628 2795.371
##          97.5% Rhat n.eff
## alpha     9.866 1.003  3000
## beta_eth[1] 0.143 1.004   850
## beta_eth[2] 0.120 1.008   380
## beta_eth[3] 0.236 1.002  3000
## beta_eth[4] 0.267 1.003  1900
## beta_height[1] 0.026 1.007   330
## beta_height[2] 0.031 1.004   560
## beta_height[3] 0.031 1.003   850
## beta_height[4] 0.034 1.005   480
## mu_beta_height 0.046 1.064  3000
## sigma      0.944 1.003   880
## sigma_height 0.104 1.037   110
## deviance   2801.508 1.006   400
##
## For each parameter, n.eff is a crude measure of effective sample size,
## and Rhat is the potential scale reduction factor (at convergence, Rhat=1).
```

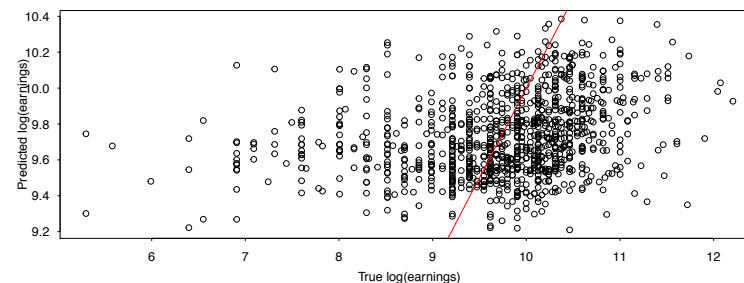
20 / 40

Interpreting the output

- ▶ The model has improved a bit, DIC now 2799.94 with 6.81 effective parameters
- ▶ The confidence intervals for the different slopes are highly different, with the whites group (ethnicity = 3) having a much clearer relationship with height, possibly due to the large sample size
- ▶ Go back to the previous classes to see plots of these effects

21 / 40

Checking the model - posterior predictive fit



22 / 40

Introducing age

- ▶ Let's fit an even more complicated model with intercepts and slopes varying by ethnicity and age group
- ▶ Age is divided up into three groups 1: 18-34, 2: 35-49, and 3: 50-64
- ▶ We want to know whether the degree to which height affects earnings for different ethnic/age group combinations

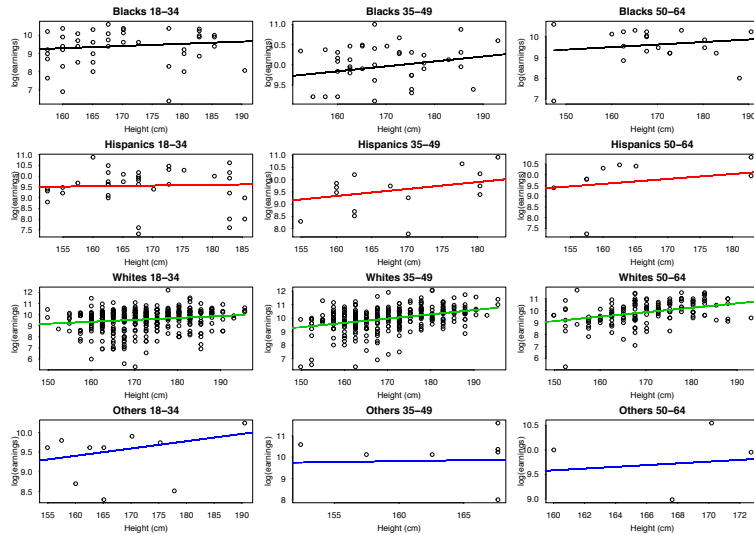
23 / 40

JAGS model

```
jags_code = '  
model{  
  # Likelihood  
  for(i in 1:N) {  
    log_earn[i] ~ dnorm(alpha[eth[i],age_grp[i]] +  
                        beta[eth[i],age_grp[i]]*(height[i] - mean(height)),  
                        sigma^-2)  
  }  
  # Priors  
  for(j in 1:N_eth) {  
    for(k in 1:N_age_grp) {  
      alpha[j,k] ~ dnorm(mu_alpha, sigma_alpha^-2)  
      beta[j,k] ~ dnorm(mu_beta, sigma_beta^-2)  
    }  
  }  
  mu_alpha ~ dnorm(0, 20^-2)  
  mu_beta ~ dnorm(0, 20^-2)  
  sigma_alpha ~ dt(0,10,1)T(0,)  
  sigma_beta ~ dt(0,10,1)T(0,)  
  sigma ~ dt(0,10,1)T(0,)  
}
```

24 / 40

Model output



25 / 40

More about this model

- ▶ So we now have varying effects - we should also plot the uncertainties in these lines (see practical)
- ▶ The DIC here is now DIC now 2737.18 with 19.42 effective parameters - a big drop!
- ▶ There is another way to make the model even more complicated . . .

26 / 40

Multivariate models

- ▶ Often we have more than one variable or parameter that changes simultaneously
- ▶ For example, we might be interested in how both earnings and social grade change in response to height. These are likely to be correlated, even after we take into account the effect of height
- ▶ Or we might be interested in how both the slope and the intercept change in a particular model
- ▶ In either case, the *multivariate normal distribution* helps us achieve this by borrowing strength across the different dimensions

27 / 40

The multivariate normal distribution

- ▶ The standard normal distribution we have met has two parameters, a mean and a standard deviation/variance
- ▶ The multivariate normal (MVN) distribution has the same, except that the mean is now a *vector* and the variance is now a *matrix* sometimes called the *covariance matrix*
- ▶ Each element of the mean represents the mean of each variable. Each diagonal element of the covariance matrix is the variance of that variable, and each off-diagonal element is the covariance between those two variables
- ▶ For the 2D MVN distribution, we write:

$$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} v_{11} & v_{12} \\ v_{12} & v_{22} \end{bmatrix} \right)$$

or, for short

$$y \sim \text{MVN}(\mu, \Sigma)$$

28 / 40

Learning about the parameters of the MVN

- ▶ You can think of the MVN distribution as borrowing strength between variables, just like the hierarchical model borrows strength across categories
- ▶ If we want to use the MVN in a Bayesian model we need to be able to put prior distributions on the parameters
- ▶ The means μ are easier as we can set independent prior distributions on them, e.g. $\mu_1 \sim N(0, 100)$
- ▶ The covariance matrix is more fiddly. Luckily there is a probability distribution called the *Wishart distribution* which works on covariance matrices. These are especially tricky as the covariances and the variances have to match appropriately

29 / 40

Returning to the earnings data

- ▶ Let's add to our complicated model with intercepts and slopes varying by ethnicity and age group
- ▶ We will include a multivariate normal prior on the intercept/slope pairs and look at the interaction between ethnicity and ages
- ▶ We are introducing a specific parameter which measures the degree of correlation between the intercept and the slope for each group. This could potentially vary between age group or ethnicity

30 / 40

Fitting a multivariate, multi-layer model

```
jags_code = '
model{
  # Likelihood
  for(i in 1:N) {
    log_earn[i] ~ dnorm(beta[eth[i],age_grp[i],1] +
                        beta[eth[i],age_grp[i],2]*(height[i] - mean(height)),
                        sigma~-2)
  }
  # Priors
  for(j in 1:N_eth) {
    for(k in 1:N_age_grp) {
      beta[j,k,1:2] ~ dnmnorm(mu_beta[1:2,1], Sigma_beta_inv)
    }
  }
  for(l in 1:2) {
    mu_beta[l,1] ~ dnorm(0, 20~-2)
  }
  Sigma_beta_inv ~ dwish(R_beta,k_beta)
  sigma ~ dt(0,10,1)T(0,)
}
```

31 / 40

Running the model

```
print(jags_run)
```

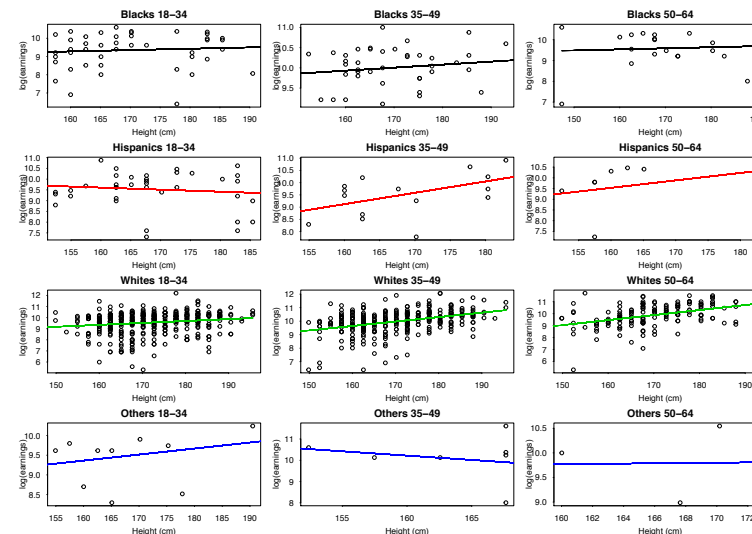
```
## Inference for Bugs model at "7", fit using jags,
## 3 chains, each with 2000 iterations (first 1000 discarded)
## n.sims = 3000 iterations saved
##          mu.vect sd.vect   2.5%    25%    50%    75%   97.5%
## beta[1,1,1]   9.354   0.131   9.097   9.269   9.355   9.441   9.616
## beta[2,1,1]   9.507   0.140   9.227   9.414   9.509   9.601   9.782
## beta[3,1,1]   9.528   0.044   9.442   9.498   9.528   9.557   9.615
## beta[4,1,1]   9.521   0.236   9.049   9.363   9.523   9.678   9.996
## beta[1,2,1]  10.008   0.134   9.744   9.916  10.004  10.100  10.265
## beta[2,2,1]   9.602   0.202   9.204   9.466   9.604   9.741   9.993
## beta[3,2,1]   9.981   0.050   9.884   9.948   9.982  10.014  10.080
## beta[4,2,1]   9.797   0.325   9.174   9.579   9.792  10.001  10.471
## beta[1,3,1]   9.590   0.176   9.240   9.476   9.589   9.707   9.937
## beta[2,3,1]   9.893   0.260   9.387   9.715   9.894  10.067  10.429
## beta[3,3,1]   9.897   0.071   9.756   9.847   9.896   9.943  10.042
## beta[4,3,1]   9.801   0.321   9.198   9.578   9.802  10.009  10.450
## beta[1,1,2]   0.007   0.013  -0.018  -0.002   0.007   0.016   0.034
## beta[2,1,2]  -0.010   0.015  -0.039  -0.019  -0.010   0.001   0.019
## beta[3,1,2]   0.017   0.005   0.009   0.014   0.017   0.021   0.026
## beta[4,1,2]   0.015   0.026  -0.037  -0.003   0.015   0.033   0.068
## beta[1,2,2]   0.007   0.014  -0.020  -0.002   0.007   0.017   0.036
## beta[2,2,2]   0.046   0.025  -0.002   0.030   0.046   0.063   0.094
## beta[3,2,2]   0.033   0.005   0.023   0.030   0.033   0.036   0.042
## beta[4,2,2]  -0.041   0.046  -0.128  -0.072  -0.041  -0.010   0.052
## beta[1,3,2]   0.005   0.017  -0.029  -0.007   0.005   0.017   0.039
## beta[2,3,2]   0.035   0.026  -0.016   0.017   0.035   0.053   0.086
## beta[3,3,2]   0.042   0.008   0.028   0.037   0.042   0.047   0.057
## beta[4,3,2]   0.003   0.080  -0.150  -0.052   0.002   0.059   0.159
## deviance      2718.778    6.659 2707.724 2714.016 2718.140 2722.704 2733.476
##          Rhat n.eff
## beta[1,1,1]  1.002  1800
## beta[2,1,1]  1.001  2500
```

32 / 40

Did the model fit the data any better?

- ▶ The DIC here is now DIC now 2740.96 with 22.18 effective parameters - about the same as before?
- ▶ We could also watch the parameters μ_{beta} and $\sigma_{\text{beta_inv}}$. The latter will tell us about the correlation structure between the intercepts and the slopes
- ▶ We could make the model even more complicated by allowing the multivariate μ_{beta} parameters to vary by age group, ethnicity, and their interaction

Plots of output



33 / 40

34 / 40

Missing and unbalanced data

- ▶ There are many definitions of what 'unbalanced' data means in statistics. Usually we mean that there are different numbers of observations in each group. Our format of writing e.g. $y_i \sim N(\alpha_{\text{eth}_i} + \beta_{\text{eth}_i} x_i, \sigma^2)$ allows us to deal with unbalanced data naturally
- ▶ Usually the smaller the sample size of the group the more uncertain the posterior distribution will be
- ▶ But what if we have some missing data? There are different types, and some need to be more carefully treated than others

Different types of missing data

- ▶ There are many different types of missing data:
 - ▶ Missing response variables
 - ▶ Missing covariates
 - ▶ Missingness that occurs completely at random
 - ▶ Missingness that occurs as a consequence of the experiment or the data
- ▶ The first three are all very easy to deal with in JAGS (less so in Stan). The last one is much harder, and not something we will go into in any detail. It requires building a separate model for the missingness process

35 / 40

36 / 40

The simple way of dealing with missing data in JAGS

- ▶ In JAGS it is absolutely trivial to deal with missingness in the response variable. You simply fill in the missing values with NA
- ▶ JAGS then treats them as parameters to be estimated. You can 'watch' them in the normal way or just ignore them. You thus have the option of getting a posterior distribution of the missing data points
- ▶ Suppose we shoved in some NA values into our data

```
dat2 = dat
dat2$earn[c(177, 763, 771)] = NA
```

37 / 40

Running the model with missingness

```
print(jags_run)

## Inference for Bugs model at "8", fit using jags,
## 3 chains, each with 2000 iterations (first 1000 discarded)
## n.sims = 3000 iterations saved
##          mu.vect sd.vect      2.5%      25%      50%      75%
## log_earn[177]  10.387   0.884   8.694   9.782  10.391  10.999
## log_earn[763]   9.454   0.872   7.799   8.843   9.451  10.061
## log_earn[771]   9.458   0.889   7.712   8.853   9.449  10.050
## deviance      2705.162   6.623 2694.066 2700.440 2704.550 2709.132
##          97.5%  Rhat n.eff
## log_earn[177]  12.135 1.001  3000
## log_earn[763]  11.151 1.002  1700
## log_earn[771]  11.222 1.002  3000
## deviance      2719.735 1.001  3000
##
## For each parameter, n.eff is a crude measure of effective sample size,
## and Rhat is the potential scale reduction factor (at convergence, Rhat=1).
##
## DIC info (using the rule, pD = var(deviance)/2)
## pD = 21.9 and DIC = 2727.1
## DIC is an estimate of expected predictive error (lower deviance is better).
```

38 / 40

More complex varieties of missing data

- ▶ If you have missing covariates or joint missing covariates/response variables, you can include these too
- ▶ The only extra issue is that you need to give JAGS a prior distribution for the missing covariate values which can make the code a bit fiddlier
- ▶ If the response variable (e.g. log earnings) exists but the covariate value is missing, then you are asking JAGS to perform an *inverse regression*
- ▶ In Stan missing data is fiddlier to incorporate as you have to separate out the parameters (i.e. missing data) from the observed data

39 / 40

Summary

- ▶ We have seen how to create some rich multi-layers models
- ▶ We have gone through quite a detailed example
- ▶ We have learnt about the multivariate normal distribution
- ▶ We have discovered how to deal with missing and unbalanced data sets

40 / 40