

# Homework #3

## Analysis of Arsenic in Rice Products

*Antoine Baldassari*

*November 17, 2015*

### 1. Standard conditionally-conjugate specification of the hierarchical model

#### 1.1 Model specification

In this model specification, the  $i^{th}$  arsenic reading of group  $j$ ,  $y_{ij}$  is normally-distributed, so that

$$y_{ij} \sim \mathcal{N}(\theta_j, \sigma^2)$$

Where  $\theta_j$  is the mean arsenic reading for the rice products indexed by  $j$ .  $\theta_j$  is normally distributed, centered at the population mean  $\mu$ , with between-group variance  $\tau^2$ :

$$\theta_j \sim \mathcal{N}(\mu, \tau^2)$$

We use conditionally-conjugate Normal and Inverse-Gamma priors on the hyperparameters:

$$\begin{aligned} 1/\sigma^2 &\sim \text{gamma}(\nu_0/2, \nu_0\sigma_0^2/2) \\ 1/\tau^2 &\sim \text{gamma}(\eta_0/2, \eta_0\tau_0^2/2) \\ \mu &\sim \text{normal}(\mu_0, \gamma_0^2) \end{aligned}$$

The full conditional distribution of the parameters can be found to be (from the book):

$$\begin{aligned} \{\theta_j \mid \sigma^2, y_{j,1}, \dots, y_{j,n}\} &\sim \mathcal{N}\left(\frac{n_j \bar{y}_j / \sigma^2 + \mu / \tau^2}{n_j / \sigma^2 + 1 / \tau^2}, [n_j / \sigma^2 + 1 / \tau^2]^{-1}\right) \\ \{\mu \mid \theta_1, \dots, \theta_m, \tau\} &\sim \mathcal{N}\left(\frac{m \bar{\theta} / \tau^2 + \mu_0 / \gamma_0^2}{m / \tau^2 + 1 / \gamma_0^2}, [m / \tau^2 + 1 / \gamma_0^2]^{-1}\right) \\ \{1/\tau^2 \mid \theta_1, \dots, \theta_m, \mu\} &\sim \mathcal{N}\left(\frac{\eta_0 + m}{2}, \frac{\eta_0 \tau_0^2 + \sum (\theta_j - \mu)^2}{2}\right) \\ \{1/\sigma^2 \mid \theta, y_1, \dots, y_n\} &\sim \mathcal{N}\left(\frac{1}{2} \left[ \nu_0 + \sum_{j=1}^m n_j \right], \frac{1}{2} \left( \nu_0 \sigma_0^2 + \sum_{j=1}^m \sum_{i=1}^{n_j} (y_{i,j} - \theta_j)^2 \right)\right) \end{aligned}$$

#### 1.2 Main analyses

We pick relatively uninformative priors, centering  $\mu$  around 1 with somewhat large within and between sample variances:  $\sigma_0^2 = 10, \nu_0 = 1, \tau_0^2 = 10, \eta_0 = 1, \gamma_0^2 = 10$ . The marginal distributions of  $\theta_1, \dots, \theta_m, \mu, \sigma^2$  and  $\tau^2$  can be obtained from the full condition distributions using a Monte-Carlo Markov-Chain algorithm,

Gibbs sampling, which we implement in R as follows::

First, we input the dataset downloaded from Sakai, modified in Stata to have numeric codes for rice products categories.

```
library(foreign)
Y <- read.dta(file="arsenicrice2.dta")
```

We set the weakly informative prior values

```
n <- nrow(Y)
nu0 <- 1; eta0 <- 1;
t20 <- 10;
mu0 <- 1;
g20 <- s20 <- var(Y$arsenic)
```

We set initial values for algorithm

```
m <- length(unique(Y$food_num)) #number of groups
n <- sv <- ybar <- rep(NA,m)
for (i in 1:m)
{
  n[i] <- sum(Y$food_num==i)
  sv[i] <- var(Y$arsenic[which(Y$food_num==i)])
  ybar[i] <- mean(Y$arsenic[which(Y$food_num==i)])
}
theta <- ybar; s2 <- mean(sv)
mu <- mean(theta); tau2 <- var(theta)
```

We create a Markov chain for each parameter by sequentially sampling from their posterior over 10,000 iterations. Elements are stored in the chain at the end of each iteration.

```
#Setup MCMC
set.seed(0808)
S <- 10000
THETA <- matrix(nrow=S, ncol=m)
OTH <- matrix(nrow=S, ncol=3)
ALL <- matrix(nrow=S, ncol=3+m)

#Run algorithm
for(i in 1:S)
{
  #Get new values for parameters
  for(j in 1:m) theta[j] <- newTheta(n[j], ybar[j], s2, tau2, mu)
  s2 <- newSigma2(m, n, nu0, s20, theta, Y)
  mu <- newMu(m, theta, tau2, g20)
  tau2 <- newTau2(m, eta0, t20, theta, mu)

  #Store in chain
  THETA[i,] <- theta
  OTH[i,] <- c(mu,s2,tau2)
  ALL[i,] <- c(theta,mu,s2,tau2)
}
```

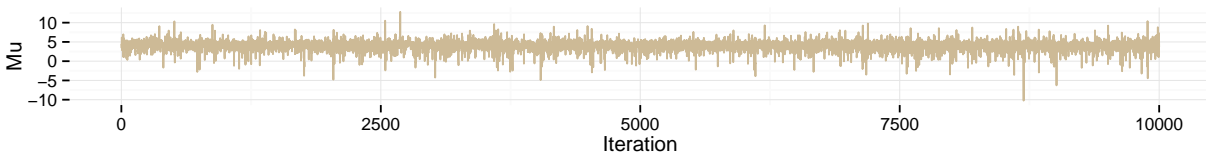
Where the functions updating the parameters follow the equations listed above:

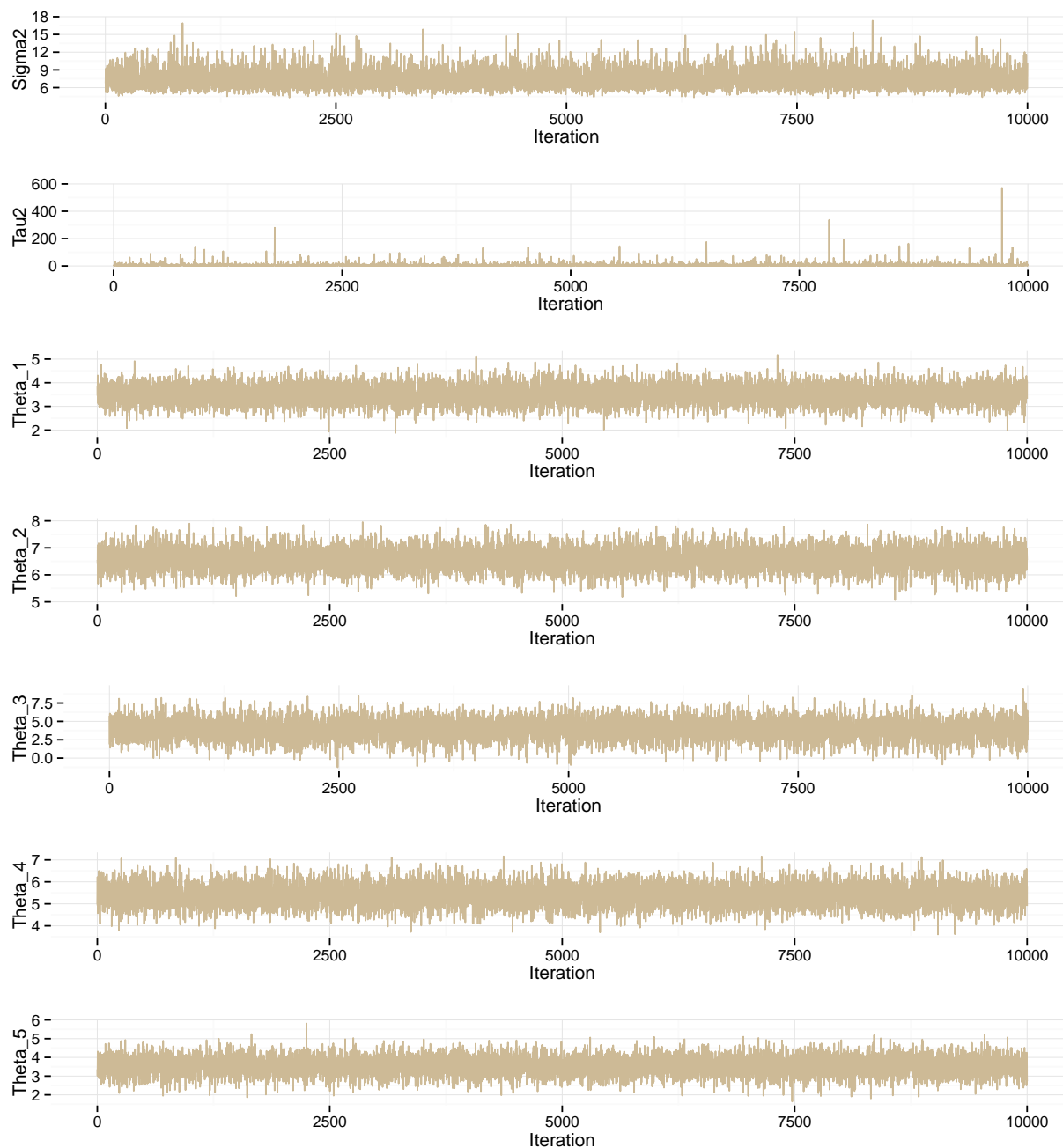
```
newTheta <- function(n, ybar, s2, tau2, mu)
{
  v = 1/(n/s2 + 1/tau2)
  e = v * (ybar*n/s2 + mu/tau2)
  new <- rnorm(1, e, sqrt(v))
  return(new)
}
newSigma2 <- function(m, n, nu0, s20, theta, Y)
{
  nun = nu0 + sum(n)
  ss <- nu0 * s20
  for(i in 1:m) ss = ss + sum((Y$arsenic[which(Y$food_num==i)] - theta[j])^2)
  sigma2 <- 1/rgamma(1, nun/2, ss/2)
  return(sigma2)
}
newMu <- function(m, theta, tau2, g20)
{
  v = 1/(m/tau2 + 1/g20)
  e = v * (m*mean(theta)/tau2 + mu0/g20)
  mu <- rnorm(1, e, v)
  return(mu)
}
newTau2 <- function(m, eta0, t20, theta, mu)
{
  etam = eta0 + m
  ss <- eta0*t20 + sum( (theta-mu) ^2 )
  tau2 <- 1/rgamma(1, etam/2, ss/2)
  return(tau2)
}
```

Before we go any further, we check that the MCMC model converged for all four statistics using ggplot2 (code used for  $\mu$  repeated for other parameters):

```
library(ggplot2)
graphdata <- data.frame(
  "Iteration"=c(1:S), "Mu"=OTH[,1], "Sigma2"=OTH[,2], "Tau2"=OTH[,3],
  "Theta_1" = THETA[,1], "Theta_2" = THETA[,2], "Theta_3" = THETA[,3],
  "Theta_4" = THETA[,4], "Theta_5" = THETA[,5])

ggplot(graphdata,aes(x=Iteration,y=Mu)) +
  theme_minimal(base_family = "") + geom_line(colour="wheat3")
```





We conclude from the graphs that convergence was achieved for all parameters.

### 1.3 Algorithm output

The estimated median values and 95% credible intervals for the parameters are as follow:

```
for(i in 1:length(ALL[1,])) print(round(unname(
  quantile(ALL[,i], probs=c(0.025, 0.5, 0.975))
),3))
```

Parameter	Credible Lower 95%	Median	Credible Upper 95%
$\theta_1$ (Basmati)	2.810	3.569	4.308
$\theta_2$ (Non-Basmati)	5.790	6.533	7.305
$\theta_3$ (Beverage)	1.869	4.231	6.448
$\theta_4$ (Cakes)	4.486	5.370	6.301
$\theta_5$ (Cereal)	2.705	3.610	4.507
$\mu$	3.193	4.697	6.182
$\sigma^2$	5.109	7.027	10.573
$\tau^2$	0.693	2.251	12.920

## 1.4 Sensitivity analyses

Evaluation of sensitivity to priors: we try three separate scenarios each tuning prior distribution of parameters:

1. Large expected  $\mu$  (Prior expectation of mad levels of arsenic)
2. Large  $\sigma^2$  and  $\nu_0$  (High variability within products)
3. Large  $\tau^2$  and  $\eta_0$  (High variability between products)

*Scenario 1:*

Parameter	Credible Lower 95%	Median	Credible Upper 95%
$\theta_1$ (Basmati)	2.706	3.488	4.265
$\theta_2^2$ (Non-Basmati)	5.902	6.671	7.460
$\theta_3$ (Beverage)	0.643	3.774	6.932
$\theta_4$ (Cakes)	4.511	5.452	6.429
$\theta_5$ (Cereal)	2.548	3.496	4.456
$\mu$	2.548	3.496	4.456
$\sigma^2$	88.793	100	111
$\tau^2$	3021.927	8427	37962

*Scenario 2:*

Parameter	Credible Lower 95%	Median	Credible Upper 95%
$\theta_1$ (Basmati)	1.237	4.396	7.069
$\theta_2^2$ (Non-Basmati)	2.811	5.401	8.559
$\theta_3$ (Beverage)	0.278	4.790	9.023
$\theta_4$ (Cakes)	1.931	4.960	8.221
$\theta_5$ (Cereal)	1.137	4.531	7.508
$\mu$	2.42	4.84	7.19
$\sigma^2$	182	215	256
$\tau^2$	0.399	1.876	22.3

*Scenario 3:*

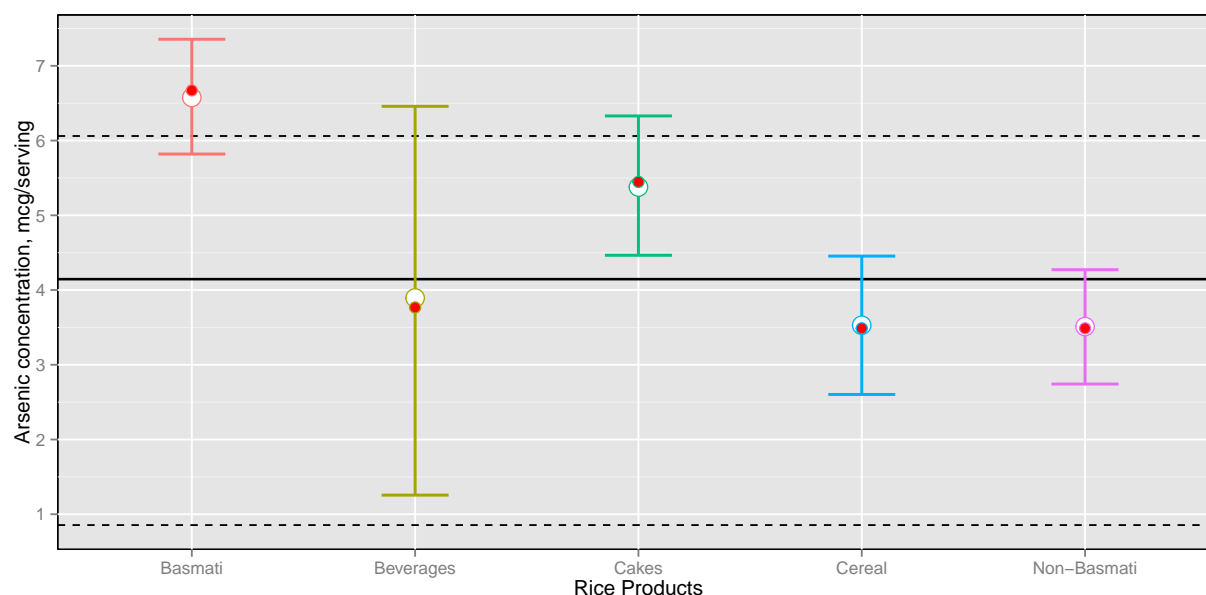
Parameter	Credible Lower 95%	Median	Credible Upper 95%
$\theta_1$ (Basmati)	2.707	3.478	4.268
$\theta_2^2$ (Non-Basmati)	5.907	6.674	7.451
$\theta_3$ (Beverage)	0.676	3.739	6.9023
$\theta_4$ (Cakes)	4.504	5.447	6.419
$\theta_5$ (Cereal)	2.547	3.494	4.452
$\mu$	-5.302	4.845	5.00
$\sigma^2$	5.197	7.327	11.37
$\tau^2$	223	288	384

We observe that excessively large prior expectations of  $\mu$  will drive up estimates of the within- and between-group variances but will have little effect on the magnitude of the estimates of within-group mean estimates (although the precision may be negatively affected for groups with relatively few observations). A large prior within-sample variance will bring posterior within-group means closer to  $\mu$ , as could be expected since the posterior estimates need to become more conservative. Increasing prior between-sample variance appears to drive up uncertainty on  $\mu$  and bring it closer to 0, without however having a notable impact on the rest of the model.

## 1.5 Results presentation

Non-Basmati rice had the highest arsenic concentration, at an estimated 6.7 mcg/serving. Rice cakes came second, at 5.4 mcg/serving, and non-Basmati and rice cereal had comparatively low amounts, slightly below 3.5 mcg/serving. There lacked data to reliably evaluate arsenic concentration in rice beverages, whose 3.8 mcg/serving estimate was particularly imprecise (95% CI=0.64, 6.93). Posterior median estimates and observed mean concentrations of arsenic are presented by product type in the following graph. Markers are  $\theta$  estimates with 95% credible interval lines; horizontal lines are the median estimate of  $\mu$  (solid) and corresponding 95% credible interval (dashed, like my hopes and dreams).

```
qmat=apply(THETA[,1:5],2,quantile,probs=c(0.025,.5,0.975))
mu_ci = quantile(OTH[,1], probs=c(0.025, 0.5, 0.975))
res <- data.frame("Rice"=c("Non-Basmati", "Basmati", "Beverages", "Cakes", "Cereal"),
                  "l95"=qmat[1,], "median"=qmat[2,], "u95"=qmat[3,], "mean"=ybar)
g <- ggplot(res, aes(x = Rice, group=Rice, colour=Rice)) +
  labs(x="Rice Products", y="Arsenic concentration, mcg/serving") +
  theme(legend.position="none", panel.background = element_rect(colour = "black")) +
  scale_y_continuous(breaks=seq(0, 7.5, 1)) +
  geom_hline(aes(yintercept=c(mu_ci[2])), size=0.7) +
  geom_hline(aes(yintercept=c(mu_ci[1])), linetype="dashed") +
  geom_hline(aes(yintercept=c(mu_ci[3])), linetype="dashed") +
  geom_errorbar(aes(ymin=l95, ymax=u95), width=.3, size=0.8) +
  geom_point(aes(y=median), fill="white", shape=21, size=5) +
  geom_point(aes(y=mean), fill="red", shape=21, size=3)
```



## 2. Parameter-expanded specification of the hierarchical model

### 2.1 Model specification

Under this model specification, instead of group means we are interested in differences between groups and the population average  $\mu$ , which is given by  $\eta_j$  for the group  $j$ , so that (under the prior belief that all groups will have equal mean):

$$\begin{aligned} y_{ij} &\sim \mathcal{N}(\mu + \xi\eta_j, \sigma_y^2) \\ \eta_i &\sim \mathcal{N}(0, \sigma_\eta^2) \end{aligned}$$

Word's on the street that well-behaved conditionally-conjugate specifications for the distributions of  $\xi$  and  $\sigma_\eta^2$  are:

$$\begin{aligned} \xi &\sim \mathcal{N}(0, 1) \\ 1/\sigma_\eta^2 &\sim \text{gamma}\left(\frac{\omega_0}{2}, \frac{\omega_0\sigma_{\eta 0}^2}{2}\right) \\ 1/\sigma_y^2 &\sim \text{gamma}\left(\frac{\nu_0}{2}, \frac{\nu_0\sigma_{y 0}^2}{2}\right) \end{aligned}$$

The prior distribution of the population mean is still  $\mu \sim \mathcal{N}(\mu_0, \gamma_0^2)$ . We set out to find full conditionals: Reparametrizing the full conditionals in exercise 1 easily yields what we need:

Note that each of the  $N$  observation in the data supports,  $\mu = y_{ij} - \xi\eta_j$  given known  $\xi$  and  $\eta_j$ 's, so that, summing over this expression and weighting it against the prior yields:

$$p(\mu, \sigma_y^2, \sigma_\eta^2, \xi, \eta_1, \dots, \eta_m | \mathbf{y}) \propto \mathcal{N}\left(\frac{\sum_j \sum_i \frac{y_{ij} - \xi\eta_j}{\sigma_y^2} + \frac{\mu_0}{\gamma_0^2}}{\frac{N}{\sigma_y^2} + \frac{1}{\gamma_0^2}}, \frac{1}{\frac{N}{\sigma_y^2} + \frac{1}{\gamma_0^2}}\right)$$

Posterior on the variances can be similarly rewritten. The sum of squares for  $\mathbf{y}$  given  $\mu, \eta_j$  and  $\xi$  is, of course,  $\sum_j \sum_i y_{ij} - \mu - \xi\eta_j$ , and the sum of squares of  $\boldsymbol{\eta}$  given its  $\mathbb{E} = 0$  is simply  $\sum_j \eta_j^2$

$$\begin{aligned} p(1/\sigma_y^2 | \mathbf{y}, \mu, \sigma_\eta^2, \xi, \eta_1, \dots, \eta_m) &\propto \text{gamma}\left(\frac{\nu_0 + N}{2}, \frac{1}{2}\left(\nu_0\sigma_{y 0}^2 + \sum_j \sum_i (y_{ij} - \mu - \xi\eta_j)^2\right)\right) \\ p(1/\sigma_\eta^2 | \mathbf{y}, \mu, \sigma_y^2, \xi, \eta_1, \dots, \eta_m) &\propto \text{gamma}\left(\frac{\omega_0 + m}{2}, \frac{1}{2}\left(\omega_0\sigma_{\eta 0}^2 + \sum_j \eta_j^2\right)\right) \end{aligned}$$

Getting full conditionals on  $\boldsymbol{\eta}$  and  $\xi$  can be likewise achieved by manipulating  $y_{ij} = \mu + \eta_j\xi$ .

$$p(\eta_j | \mathbf{y}, \mu, \sigma_y^2, \sigma_\eta^2, \xi) \propto \mathcal{N} \left( \frac{\xi \sum_i y_{ij} - \mu}{\frac{\sigma_y^2}{\sigma_\eta^2} + \frac{1}{\sigma_\eta^2}}, \frac{1}{\frac{n_j \xi^2}{\sigma_y^2} + \frac{1}{\sigma_\eta^2}} \right)$$

and

$$p(\xi | \mathbf{y}, \mu, \sigma_y^2, \sigma_\eta^2, \eta_1, \dots, \eta_m) \propto \mathcal{N} \left( \frac{\sum_j \eta_j \sum_i (y_{ij} - \mu)}{\frac{\sigma_y^2}{\sigma_\eta^2} + 1}, \frac{1}{\sum_j \frac{n_j \eta_j^2}{\sigma_y^2} + 1} \right)$$

## 2.2 Analyses

Similarly to exercise 1, we pick the priors  $\sigma_{y0}^2 = 10, \nu_0 = 1, \omega_0 = 1, \sigma_{0\eta}^2 = 10$  and  $\gamma_0^2 = 10$ , and proceed with Gibbs sampling:

We read the data

```
library(foreign)
library(hdrcde)
```

```
## Loading required package: mvtnorm
## hdrcde 3.1 loaded
```

```
data <- read.dta(file="arsenicrice2.dta")
Y <- read.dta(file="arsenicrice2.dta")
```

We set prior values

```
n <- nrow(Y)
nu0 <- 1; omega0 <- 1;
s2_eta0 <- 10; s2_y0 <- 10
mu0 <- mean(Y$arsenic);
g20 <- var(Y$arsenic)
```

We setup the MCMC

```
#Setup starting values
m <- length(unique(Y$food_num)) #number of groups
n <- sv_y <- ybar <- rep(NA,m) #create empty vectors for group descriptions
for (i in 1:m)
{
  n[i] <- sum(Y$food_num==i)
  sv_y[i] <- var(Y$arsenic[which(Y$food_num==i)])
  ybar[i] <- mean(Y$arsenic[which(Y$food_num==i)])
}
eta <- ybar - mean(Y$arsenic); s2_y <- mean(sv_y)
```



```

mu <- mean(Y$arsenic); s2_eta <- var(eta)
xi <- 0

#Setup MCMC
set.seed(0808)
S <- 10000
THETA <- matrix(nrow=S, ncol=m)
RES <- matrix(nrow=S, ncol=4)
ALL <- matrix(nrow=S, ncol=4+m)

#Setup MCMC
set.seed(0808)
S <- 10000
ETA <- matrix(nrow=S, ncol=m)
RES <- matrix(nrow=S, ncol=4)
ALL <- matrix(nrow=S, ncol=4+m)

```

Our updating functions correspond to the full conditionals derived above:

```

#FUNCTIONS
newS2eta <- function(m, s2_eta0, eta, omega0)
{
  ss <- s2_eta0*omega0 + sum( eta^2 )
  s2_eta <- 1/rgamma(1, shape=((omega0 + m)/2), scale=(ss/2))
  return(s2_eta)
}
newS2y <- function(nu0, n, s2_y0, y, m, eta, xi, mu)
{
  nun = nu0 + sum(n)
  ss = sum((y - mu - xi*rep(eta, times=n))^2) + nu0*s2_y0
  s2_y <- 1/rgamma(1, shape=(nun/2), scale=(ss/2))
  return(s2_y)
}
newMu <- function(y, xi, eta, s2_y, g20, mu0, n)
{
  v = ( sum(n)/s2_y + 1/g20 )
  ss = sum(y-rep(eta,times=n)*xi)/s2_y + mu0/g20
  e = (ss/s2_y + mu0/g20)
  mu = rnorm(1, ss/v, sqrt(1/v))
  return(mu)
}
newEta <- function(xi, mu, ybar, n, s2_y, s2_eta)
{
  v = 1/(n*xi^2/s2_y + 1/s2_eta)
  e = v*xi*(n*ybar - mu*n)/s2_y
  eta = rnorm(m, e, sqrt(v))
  return(eta)
}
newXi <- function(eta, m, y, n, mu, s2_y)
{
  v = (sum(n*eta^2)/s2_y + 1)
  e = sum((y-mu)*rep(eta,times=n))/s2_y
  xi = rnorm(1, e/v, sqrt(1/v))
}

```

```

    return(xi)
}

```

We run the MCMC algorithm:

```

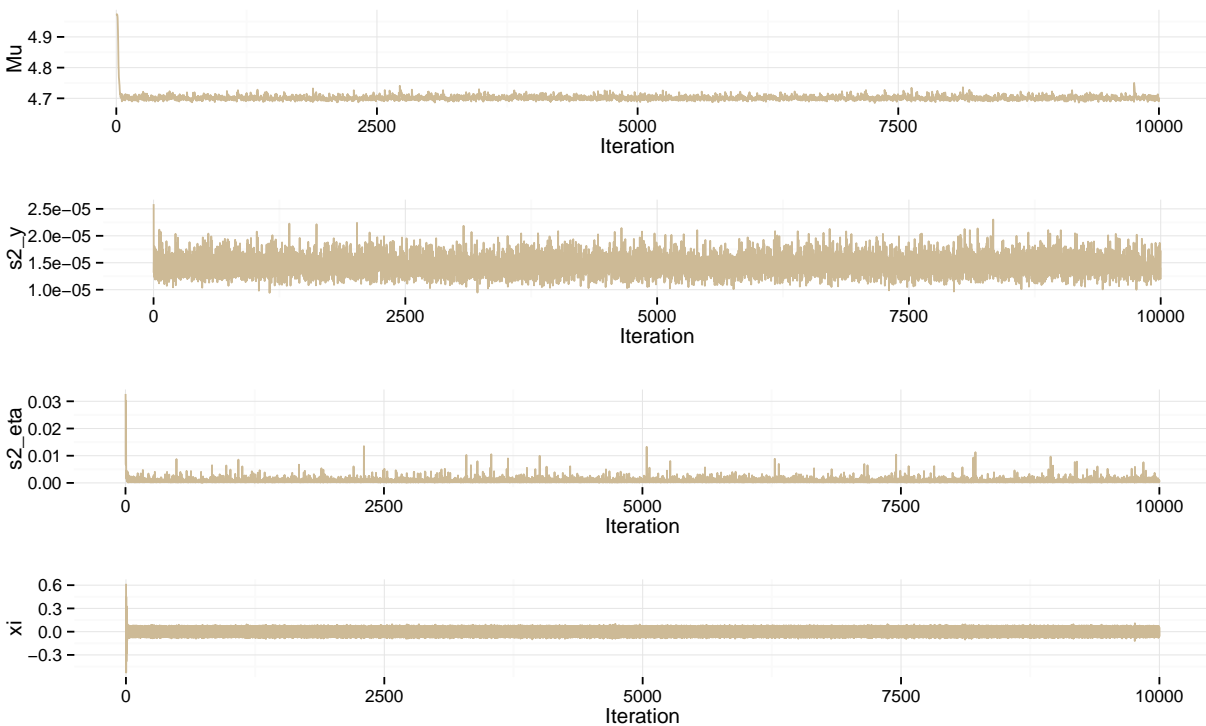
#RUN MCMC
for(i in 1:S)
{
  #Get new values for parameters
  eta <- newEta(xi, mu, ybar, n, s2_y, s2_eta)
  mu <- newMu(Y$arsenic, xi, eta, s2_y, g20, mu0, n)

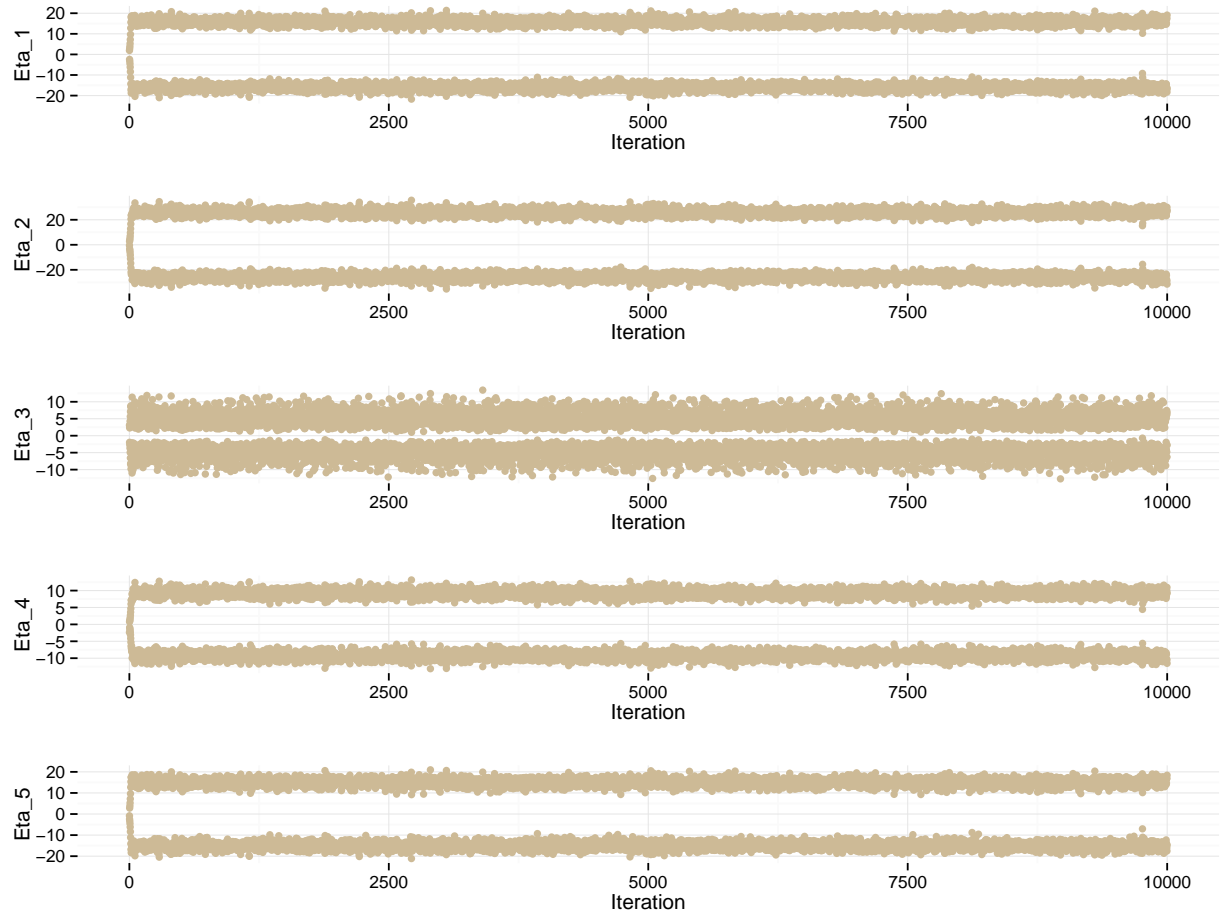
  s2_eta <- newS2eta(m, s2_eta0, eta, omega0)
  s2_y <- newS2y(nu0, n, s2_y0, Y$arsenic, m, eta, xi, mu)

  xi <- newX(eta, m, Y$arsenic, n, mu, s2_y)
  #Store in chain
  ETA[i,] <- eta
  RES[i,] <- c(mu, s2_y, s2_eta, xi)
  ALL[i,] <- c(eta, mu, xi, s2_eta, s2_y)
}

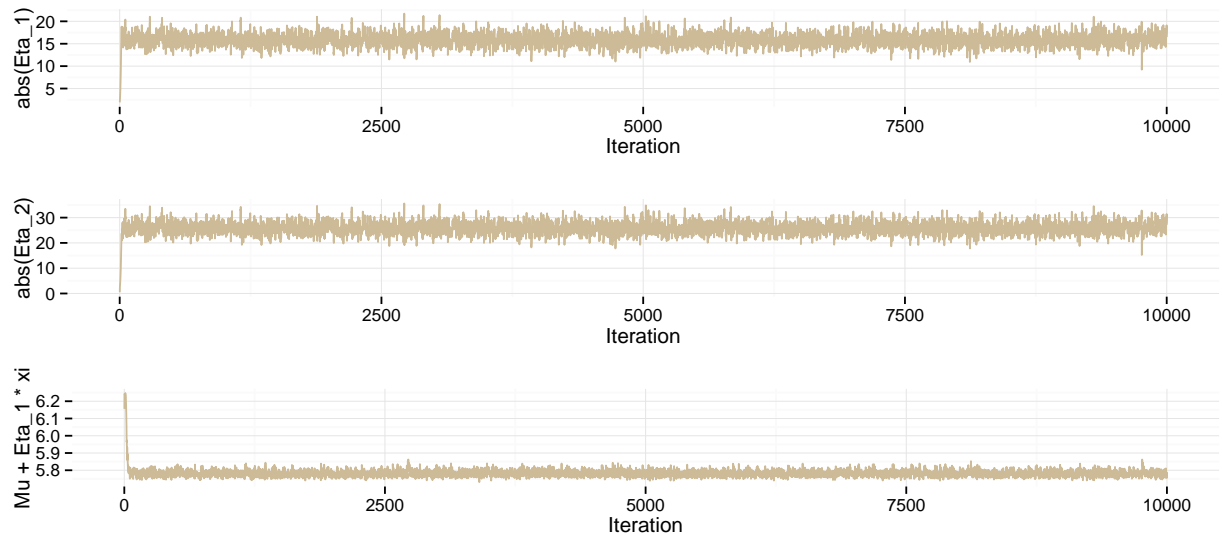
```

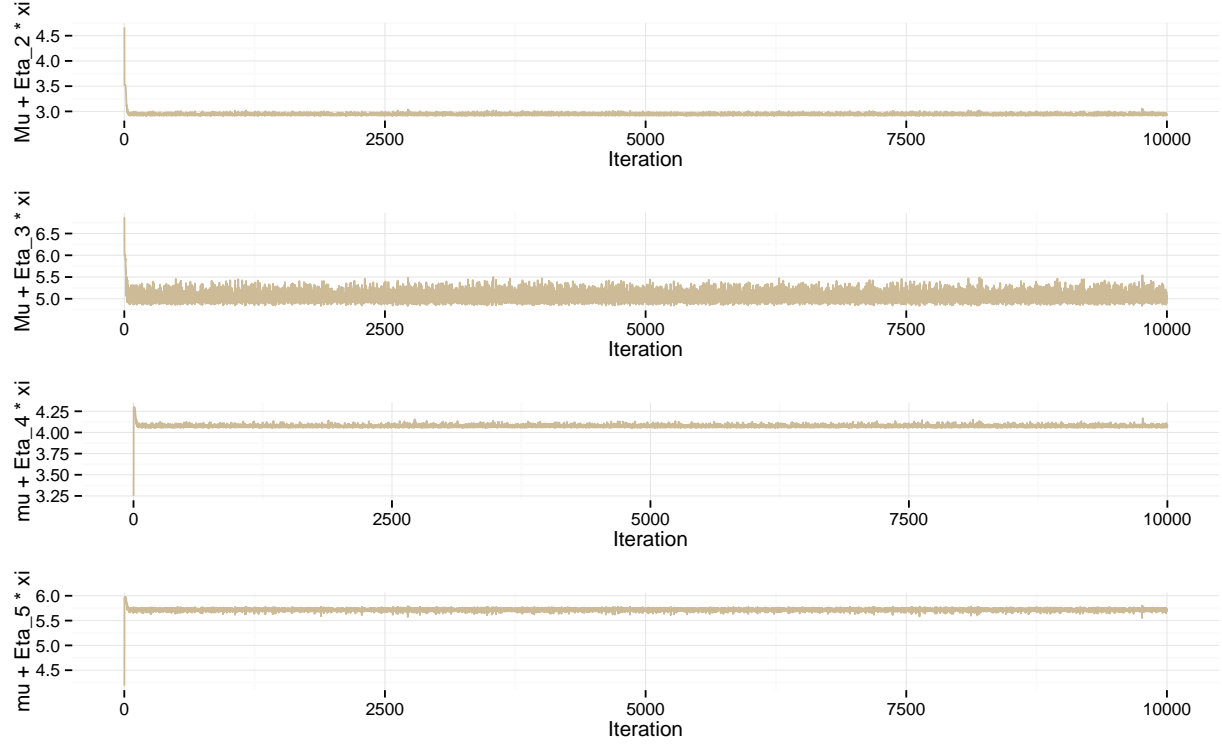
Again, before we get too psyched about results, we check MCMC convergence criteria (code not shown, see above):





Obviously, something is wrong with the  $\eta$  parameters, which could make sense since when the estimate of  $\xi$  crosses zero, the  $\eta$  parameters get updated on the other side of 0 as well. Taking the absolute value of  $\xi$  reassures us that the parameter space that is actually searched isn't terribad. We also check the convergence of  $\mu + \xi\eta$ , which is the mean concentration of arsenic in each group, and ultimately interests us. Reproduced below are the graphs for  $|\eta_1|$ ,  $|\eta_2|$ , and  $\theta$





We are fairly satisfied with the outlook on the convergence of our  $\theta$ 's.

## 2.3 Algorithm output

We provide the median estimate and 95% credible interval for the  $\theta$  parameters in this expanded hierarchical specification model:

Parameter	Credible Lower 95%	Median	Credible Upper 95%
$\theta_1$ (Basmati)	2.810	3.569	4.308
$\theta_2$ (Non-Basmati)	5.790	6.533	7.305
$\theta_3$ (Beverage)	1.869	4.231	6.448
$\theta_4$ (Cakes)	4.486	5.370	6.301
$\theta_5$ (Cereal)	2.705	3.610	4.507
$\mu$	3.193	4.697	6.182
$\sigma^2$	5.109	7.027	10.573
$\tau^2$	0.693	2.251	12.920

## 2.4 Sensitivity analyses

### **3. Conditionally-conjugate specification of the hierarchical model with group-specific variances**

#### **3.1 Model specification**

#### **3.2 Analyses**

#### **3.3 Algorithm output**

#### **3.4 Sensitivity analyses**