

Problem 1:

- (a)  $\text{Var}(y_{ij} | u, v^2)$  is bigger, because it includes both within and between group variability, while  $\text{Var}(y_{ij} | \theta_j, \sigma^2)$  only includes within group variability.

- (b) if  $\theta_j$  and  $\sigma^2$  is known,

Since  $y_{1j}, y_{2j} | \theta_j, \sigma^2 \sim N(\theta_j, \sigma^2)$

$\text{Cov}(y_{1j}, y_{2j} | \theta_j, \sigma^2) = 0$ , i.e., knowing  $y_{1j}$  gives no information about  $y_{2j}$  if  $\theta_j$  is known.

if  $\theta_j$  is unknown,

knowing about  $y_{1j}$  gives information about  $\theta_j$ , thus gives information about  $y_{2j}$ .

if  $y_{1j}$  is large, we expect  $\theta_j$  is large, and  $y_{2j}$  is large.

$\text{Cov}(y_{1j}, y_{2j} | u, v^2) > 0$

- (c)  $\text{① } \text{Var}(y_{ij} | \theta_j, \sigma^2) = \sigma^2$

$$\text{② } \text{Var}(\bar{y}_{ij} | \theta_j, \sigma^2) = \text{Var}\left(\frac{\sum_{i=1}^{n_j} y_{ij}}{n_j} | \theta_j, \sigma^2\right) = \frac{\sigma^2}{n_j}$$

$$\begin{aligned} \text{③ } \text{Cov}(y_{1j}, y_{2j} | \theta_j, \sigma^2) &= E(y_{1j} y_{2j} | \theta_j, \sigma^2) - E(y_{1j} | \theta_j, \sigma^2) E(y_{2j} | \theta_j, \sigma^2) \\ &= 0 \end{aligned}$$

$$\text{④ } \text{Var}(y_{ij} | u, v^2) = \text{Var}(E(y_{ij} | \theta_j, \sigma^2) | u, v^2) + E(\text{Var}(y_{ij} | \theta_j, \sigma^2) | u, v^2)$$

$$= \text{Var}(\theta_j | u, v^2) + E(\sigma^2 | u, v^2)$$

$$= \tau^2 + \sigma^2$$

$$\text{⑤ } \text{Var}(\bar{y}_{ij} | u, v^2) = \text{Var}(E(y_{ij} | \theta_j, \sigma^2) | u, v^2) + E(\text{Var}(y_{ij} | \theta_j, \sigma^2) | u, v^2)$$

$$= \text{Var}(\theta_j | u, v^2) + E\left(\frac{\sigma^2}{n_j} | u, v^2\right)$$

$$= \tau^2 + \frac{\sigma^2}{n_j}$$

$$\text{⑥ } \text{Cov}(y_{1j}, y_{2j} | u, v^2) = E(\text{Cov}(y_{1j}, y_{2j} | \theta_j, \sigma^2) | u, v^2) + \text{Cov}[E(y_{1j} | \theta_j, \sigma^2), E(y_{2j} | \theta_j, \sigma^2) | u, v^2]$$

$$= \text{Cov}(\theta_j, \theta_j | u, v^2)$$

$$= \text{Var}(\theta_j | u, v^2)$$

$$= \tau^2$$

$$\text{Var}(y_{ij} | \theta_j, \sigma^2) < \text{Var}(y_{ij} | u, v^2)$$

$$\text{Cov}(y_{1j}, y_{2j} | \theta_j, \sigma^2) = 0$$

$$\text{Cov}(y_{1j}, y_{2j} | u, v^2) > 0$$

which align with the results in (a) and (b)

$$(d) p(u | \theta_1, \dots, \theta_m, \sigma^2, \tau^2, y_1, \dots, y_m)$$

$$= \frac{p(u, \theta_1, \dots, \theta_m, \sigma^2, \tau^2, y_1, \dots, y_m)}{\int p(u, \theta_1, \dots, \theta_m, \sigma^2, \tau^2, y_1, \dots, y_m) du}$$

$$= \frac{p(y_1, \dots, y_m | \theta_1, \dots, \theta_m, \sigma^2) p(\theta_1, \dots, \theta_m | u, v^2) p(u) p(\sigma^2, \tau^2)}{p(y_1, \dots, y_m | \theta_1, \dots, \theta_m, \sigma^2) p(\sigma^2, \tau^2) \int p(\theta_1, \dots, \theta_m | u, v^2) p(u) du}$$

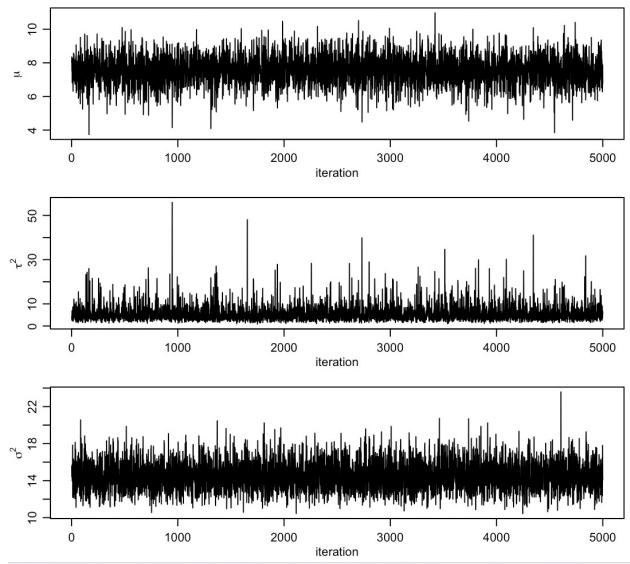
$$= \frac{p(\theta_1, \dots, \theta_m | u, v^2) p(u)}{\int p(\theta_1, \dots, \theta_m | u, v^2) p(u) du}$$

$$= p(u | \theta_1, \dots, \theta_m, \tau^2)$$

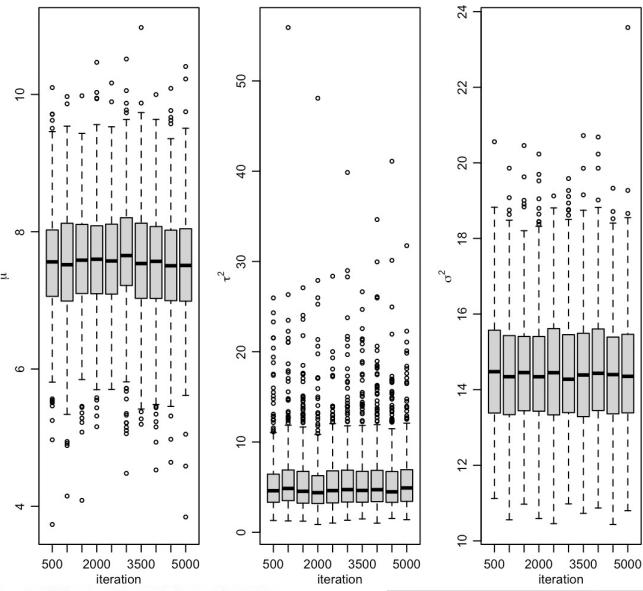
This means  $\mu$  does not depend on  $y_1, \dots, y_n$  and  $\sigma^2$ , i.e.,  $\mu$  does not depend on the observed data, if  $\theta_1, \dots, \theta_m$  is known.

Problem 2:

(a) trace plot:



stationary plot:



> effectiveSize(MU)

var1  
4147.251

effective sample size of  $\mu$

> effectiveSize(TAU\_SQ)

var1

effective sample size of  $\tau^2$

3729.928

> effectiveSize(SIGMA\_SQ)

var1

effective sample size of  $\sigma^2$

4713.354

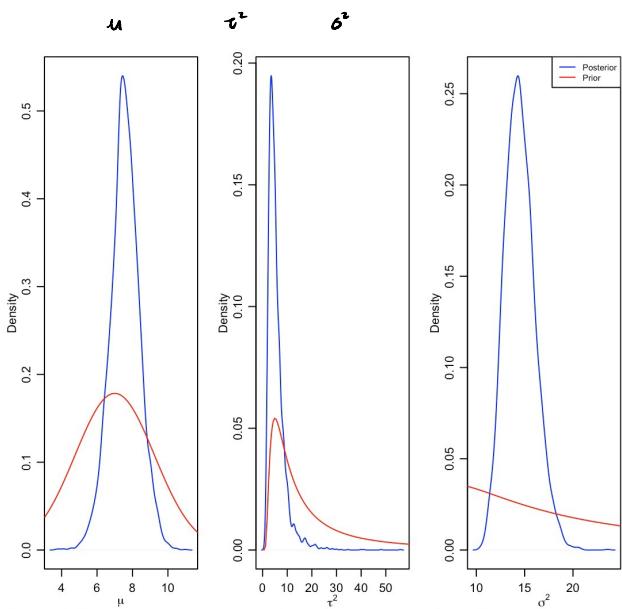
From the stationary plot and trace plot, we can see our MCMC samples reaches stationary. effective sample sizes are all  $> 1000$ , so the number of iterations is large enough.

(b) posterior means of  $\mu, \tau^2, \sigma^2$

```
> apply(M, MARGIN=2, FUN=mean)
[1] 7.563386 5.581711 14.493790
```

posterior confidence regions of  $\mu$ ,  $\tau^2$ ,  $\sigma^2$  respectively.

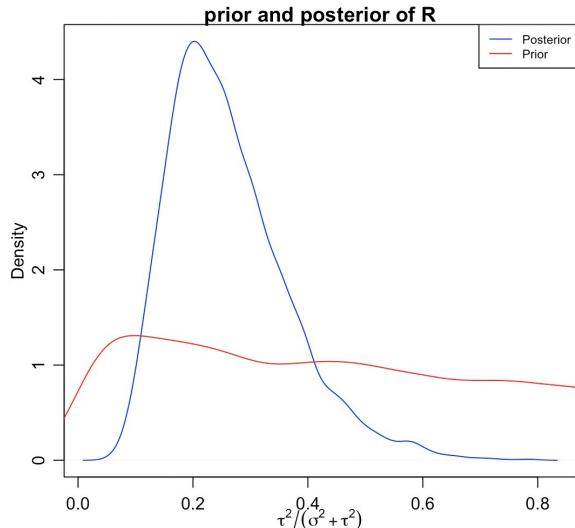
```
> apply(M, MARGIN = 2, FUN = quantile, probs = c(0.025, 0.975))
[,1]      [,2]      [,3]
2.5%  5.877648 1.83903 11.76387
97.5% 9.200664 15.28747 17.82563
```



remark: the posteriors are more centered than prior distribution.

So after observing the data, we are more confident about the true value of  $\mu$ ,  $\tau^2$  and  $\sigma^2$ .  $\mu$  is around 7.5,  $\tau^2$  is around 5.5, and  $\sigma^2$  is around 14.5.

(c)



posterior mean of R  
 $> \text{mean}(R.\text{posterior})$   
[1] 0.2618374

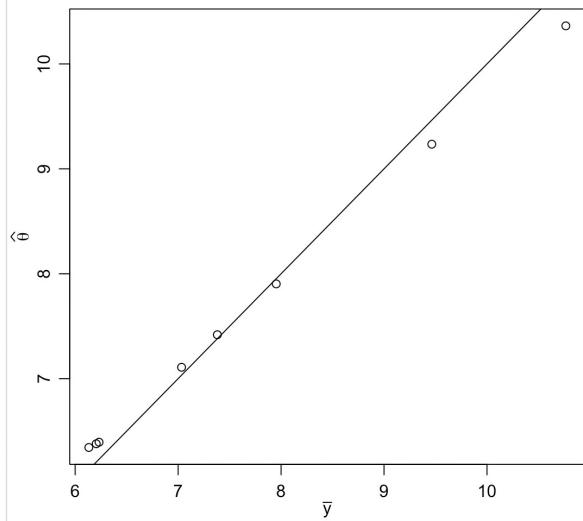
$R = \frac{\tau^2}{\sigma^2 + \tau^2}$  represents the proportion of the between school variability among total variability of  $y_{ij}$ .

After observing the data, the posterior mean of  $R$  is 0.26, which means the within group variability contributes most of the total variability.

(d)  $p(\theta_1 < \theta_6 | y)$      $> \text{mean}(\text{THETA}[7] < \text{THETA}[6])$   
[1] 0.5208

$p(\theta_1 < \theta_6 | y)$      $> \text{mean}(\text{minimum} == 7)$   
[1] 0.3208

(e)



relationship of  $\bar{y}_j$  and  $\hat{\theta}_j$  does not follow  $y=x$ .

① when  $\bar{y}_j$  is away from  $u$ , there will be shrinkage effect:

$\hat{\theta}_j$  moves close to  $u$ , difference between  $\bar{y}_j$  and  $\hat{\theta}_j$  will be larger if  $|\bar{y}_j - u|$  is larger.

② when  $\bar{y}_j$  is close to  $u$ ,  $\bar{y}_j$  is also close to  $\hat{\theta}_j$ .

```
> sum/sum(n)      →  $\bar{y}$ 
[1] 7.691278
> mean(MU)
[1] 7.563386      →  $u$ 
```

### Problem 3:

(a) posterior mean of  $\beta$ :

```
> beta.posterior
[1] 0.277807220 0.002181378 0.532223831 1.438391578 -0.764108880 -0.068558382 0.134361426 -0.065762314
[9] 0.108240773 -0.266088740 0.358137273 0.229835939 0.708270476 -0.279304157 -0.058860836
```

95% CI of  $\beta$ :

```
> apply(BETA, 2, quantile, probs = c(0.025, 0.975))
[,1]      [,2]      [,3]      [,4]      [,5]      [,6]      [,7]      [,8]      [,9]
2.5% 0.02763848 -0.3274552 0.2068963 -0.06838007 -2.297244 -0.3407333 -0.1424860 -0.2882754 -0.2198412
97.5% 0.52641303 0.3392396 0.8583802 2.89560840 0.808366 0.2153559 0.4140429 0.1641163 0.4265987
[,10]     [,11]     [,12]     [,13]     [,14]     [,15]
2.5% -0.61742527 0.02330779 -0.2339990 0.2991445 -0.5242016 -0.2950214
97.5% 0.09618591 0.67841679 0.6969059 1.1266416 -0.0372050 0.1901738
```

M, Ed, Uz, Ineq, Prob are strongly predictive of crime rates.

because their 95% credible intervals do not contain 0.

M, Ed, Uz, Ineq have significant positive relationship with crime rates.

and Prob has significant negative relationship with crime rates.

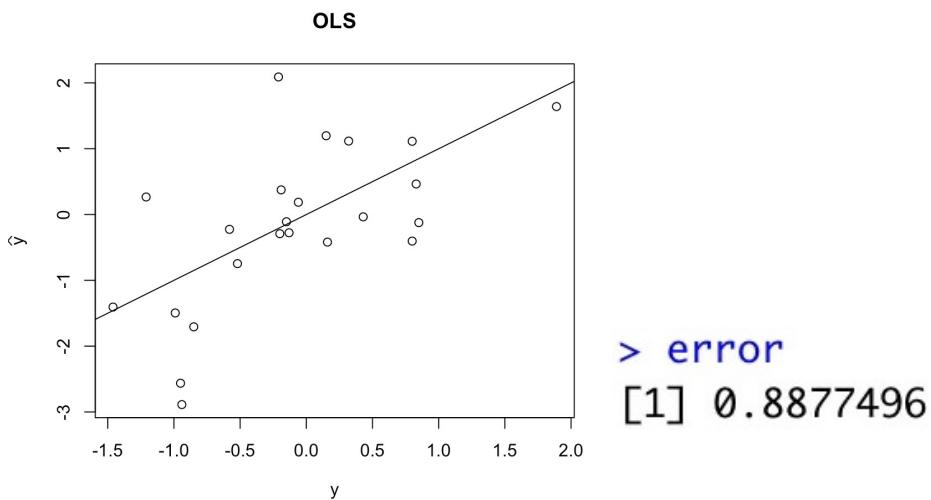
Other variables do not have significant relationship with crime rates.

least square estimates for  $\beta$ :

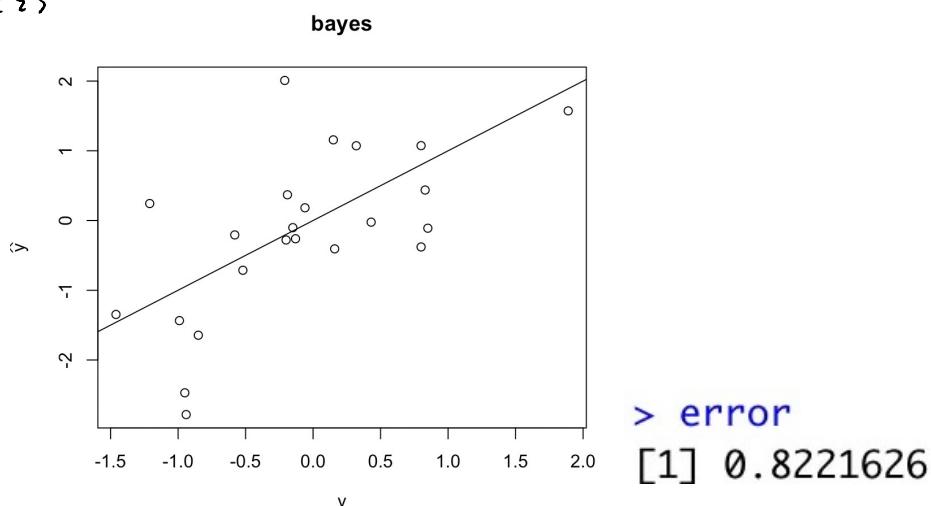
```
> beta.ols  
[1]  
M      0.2865177028  
So     -0.0001179958  
Ed     0.5445161778  
Po1    1.4716146465  
Po2    -0.7817757455  
LF     -0.0659672893  
M.F    0.1313002714  
Pop   -0.0702910179  
NW    0.1090590127  
U1    -0.2705407273  
U2    0.3687335028  
GDP   0.2380580097  
Ineq  0.7262918200  
Prob  -0.2852262729  
Time -0.0615771841
```

least square estimates are similar to marginal posterior means of  $\beta$ .

(b) <1>



<2>



$\hat{\beta}_{\text{Bayes}}$  has slightly better accuracy for the test data than  $\hat{\beta}_{\text{OLS}}$

$\hat{\beta}_{\text{OLS}}$  works best for training data, but may not work well for test data.

(c) We run 50 times, i.e., sample 50  $\{Y_{it}, X_{it}\}$ 's and  $\{Y_{it}, X_{it}\}$ 's.

> mean>Error\_OLS → prediction error of  $\hat{\beta}_{OLS}$

[1] 0.8174934

> mean>Error\_Bayes → prediction error of  $\hat{\beta}_{Bayes}$ .

[1] 0.783311

> mean>Error\_Bayes < Error\_OLS  
[1] 0.9

Average prediction error is smaller for Bayes estimator than OLS estimator.  
for 90% of the times among 50 iterations, prediction error of  $\hat{\beta}_{Bayes}$  is smaller.

## Problem 4

$$(a) \theta_i = \Pr(Y_i=1 | \alpha, \beta, x_i)$$

$$\log\left(\frac{\theta_i}{1-\theta_i}\right) = \alpha + \beta x_i$$

$$\therefore \theta_i(1 + e^{\alpha + \beta x_i}) = e^{\alpha + \beta x_i}$$

$$\theta_i = \frac{e^{\alpha + \beta x_i}}{1 + e^{\alpha + \beta x_i}}$$

$$\begin{aligned} \therefore \prod_{i=1}^n p(Y_i | \alpha, \beta, x_i) &= \prod_{i=1}^n \left( \frac{e^{\alpha + \beta x_i}}{1 + e^{\alpha + \beta x_i}} \right)^{y_i} \left( \frac{1}{1 + e^{\alpha + \beta x_i}} \right)^{1-y_i} \\ &= \prod_{i=1}^n \frac{e^{(\alpha + \beta x_i) y_i}}{1 + e^{\alpha + \beta x_i}} \end{aligned}$$

(b) we want  $0 \leq \theta_i \leq 1$  for  $10 \leq x_i \leq 15$

$$\text{i.e., } 0 \leq \frac{e^{\alpha + \beta x_i}}{1 + e^{\alpha + \beta x_i}} \leq 1 \text{ for } 0 \leq x_i \leq 15$$

we want prior be uninformative,

so we let  $10^{-3} \leq \theta_i \leq 1 - 10^{-3}$  be the prior,

$$\therefore \text{approximately, } -5 \leq \log \frac{\theta_i}{1-\theta_i} \leq 5$$

$$\therefore -5 \leq \alpha + \beta x_i \leq 5$$

$$\text{when } \beta = 0, \quad -5 \leq \alpha \leq 5$$

$$\text{when } x_i = 10 \text{ be minimum, } -\frac{1}{2} \leq \beta \leq \frac{1}{2}$$

$$\therefore \beta \sim N(0, \frac{1}{4}) \text{ be prior of } \beta$$

$$\text{when } \beta = 0, \quad -5 \leq \alpha \leq 5$$

$$\therefore \alpha \sim N(0, 25) \text{ be prior of } \alpha$$

(c) I ran the Metropolis algorithm 50000 times.

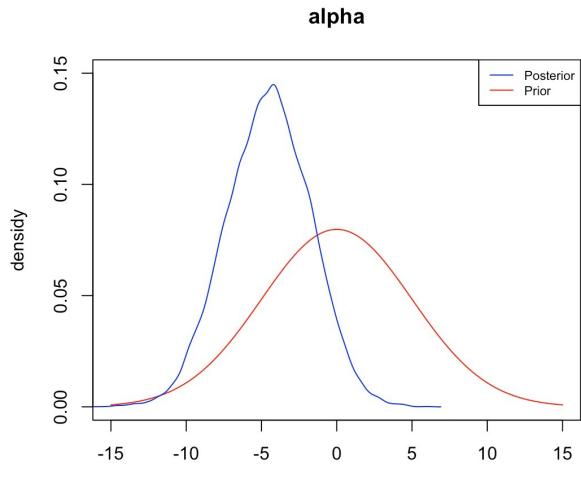
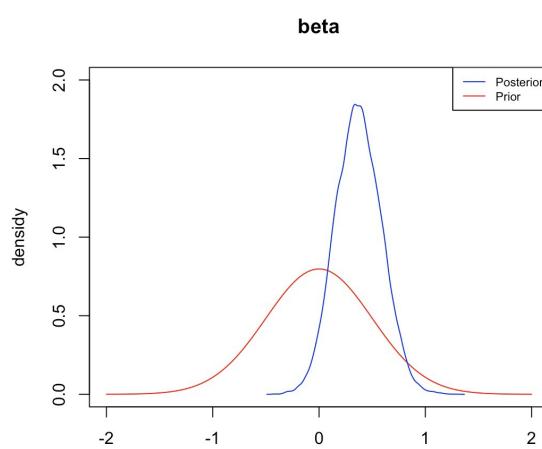
the effective sample size is:

> apply(THETA, 2, effectiveSize)

[1] 1359.971 1346.402

$$\begin{array}{ccc} \uparrow & & \uparrow \\ \alpha & & \beta \end{array}$$

(d)

posterior and prior for  $\alpha$ :posterior and prior for  $\beta$ .

posterior of  $\alpha$  and  $\beta$  are both more centered than their own priors.

(e)

