

## STAT 289 Homework 3

Wenjing Xu

For Questions 1-7, you will need the SAT coaching dataset given in Table 1 below. This dataset comes from a randomized experiment to test the effects of 8 different SAT coaching programs. For each program  $i$ , an observed effect  $y_i$  was measured as well as a standard error  $\sigma_i$  of the treatment effect. We want to analyze these treatments under a hierarchical model where

$$\begin{aligned} y_i &\sim \text{Normal}(\theta_i, \sigma_i^2) \\ \theta_i &\sim \text{Normal}(\mu, \tau^2) \end{aligned}$$

We will assume that  $\boldsymbol{\sigma}^2 = (\sigma_1^2, \dots, \sigma_n^2)$  are all known without error. We are interested in posterior inference for each  $\theta_i$ , which is the underlying true treatment effect for program  $i$ , as well as the common treatment mean  $\mu$  and variance  $\tau^2$ .

1. Verify that using a flat prior of  $p(\mu, \tau) \propto 1$  corresponds to a prior for  $(\mu, \tau^2)$  of  $p(\mu, \tau^2) \propto \tau^{-1}$ .

With this prior, write out the full posterior distribution of the unknown parameters  $\mu$ ,  $\tau^2$  and  $\boldsymbol{\theta} = (\theta_1, \dots, \theta_n)$ .

Proof: Since  $p(\mu, \tau) \propto 1$ , we can denote  $t = \tau^2$ , then  $\tau = \sqrt{t} \Rightarrow \frac{d\tau}{dt} = \frac{1}{2\sqrt{t}} \Rightarrow p(\mu, t) \propto \frac{1}{\sqrt{t}}$ , since  $t = \tau^2$ , then we can get  $p(\mu, \tau^2) \propto \frac{1}{\tau}$

Denote  $\boldsymbol{\theta} = (\theta_1, \theta_2, \dots, \theta_n)$ ,  $\mathbf{y} = (y_1, y_2, \dots, y_n)$ ,  $\boldsymbol{\sigma}^2 = (\sigma_1^2, \sigma_2^2, \dots, \sigma_n^2)$ , then the full posterior distribution of  $\mu, \tau^2, \boldsymbol{\theta}$  is:

$$\begin{aligned} p(\mu, \tau^2 | \mathbf{y}, \boldsymbol{\theta}^2) &\propto p(\mu, \tau^2) p(\boldsymbol{\theta} | \mu, \tau^2) p(\mathbf{y} | \boldsymbol{\theta}, \boldsymbol{\sigma}^2) \\ &\propto \frac{1}{\tau} \prod_{i=1}^n [N(\theta_i | \mu, \tau^2) N(y_i | \theta_i, \sigma_i^2)] \end{aligned}$$

2. The first approach is based on the realization that we can actually integrate out each  $\theta_i$  from the above distribution, giving us the marginal distribution of our treatment effects  $y_i$ :

$$y_i \sim \text{Normal}(\mu, \sigma_i^2 + \tau^2)$$

Write out the marginal posterior distribution of  $p(\mu, \tau^2 | \mathbf{y}, \boldsymbol{\sigma}^2)$ , and evaluate this posterior distribution over a grid of values of  $\mu$  and  $\tau^2$ . Use a grid of (0,10) for  $\tau^2$ .

Solution: The marginal distribution of  $y_i$  is  $N(\mu, \sigma_i^2 + \tau^2)$ , then the marginal posterior:

$$\begin{aligned} p(\mu, \tau^2 | \mathbf{y}, \boldsymbol{\sigma}^2) &\propto p(\mu, \tau^2) p(\mathbf{y} | \mu, \tau^2, \boldsymbol{\sigma}^2) \\ &\propto \frac{1}{\tau} \prod_{i=1}^n N(y_i | \mu, \sigma_i^2 + \tau^2) \\ &\Rightarrow p(\mu | \tau^2, \mathbf{y}, \boldsymbol{\sigma}^2) \propto e^{-\sum \frac{(y_i - \mu)^2}{2(\sigma_i^2 + \tau^2)}} \\ &\Rightarrow p(\mu | \tau^2, \mathbf{y}, \boldsymbol{\sigma}^2) = N(\mu | \beta, V^2), \end{aligned}$$

where  $V^2 = (\sum \frac{1}{\sigma_i^2 + \tau^2})^{-1}$ ,  $\beta = (\sum \frac{y_i}{\sigma_i^2 + \tau^2})V^2$ .

$$\Rightarrow p(\tau^2 | \mathbf{y}, \boldsymbol{\sigma}^2) = \frac{p(\mu, \tau^2 | \mathbf{y}, \boldsymbol{\sigma}^2)}{p(\mu | \tau^2, \mathbf{y}, \boldsymbol{\sigma}^2)} \propto \frac{\frac{1}{\tau} \prod_{i=1}^n N(y_i | \mu, \sigma_i^2 + \tau^2)}{N(\mu | \beta, V^2)}$$

3. Use the grid sampling method to get 1000 samples from  $p(\mu, \tau^2 | \mathbf{y}, \boldsymbol{\sigma}^2)$ . Calculate the mean, median and 95% posterior interval for  $\mu$  and  $\tau^2$ .

Solution: 1000 samples were got from  $p(\mu, \tau^2 | \mathbf{y}, \boldsymbol{\sigma}^2)$  by using the grid sampling method.

	mean	median	95% posterior interval
$\mu$	0.01546715	0.01252507	[-0.44199333, 0.46522200]
$\tau^2$	9.609489	9.709739	[8.668301, 9.989991]

4. Write out the distribution  $p(\theta_i | \mu, \tau^2, \mathbf{y}, \boldsymbol{\sigma}^2)$ . Use the 1000 samples of  $\mu$  and  $\tau^2$  to draw 1000 samples of each  $\theta_i$ . Calculate the mean of each  $\theta_i$  and compare to the observed treatment effects  $y_i$ .

Solution: Since  $p(\mu, \tau^2, \boldsymbol{\theta} | \mathbf{y}, \boldsymbol{\sigma}^2) \propto \frac{1}{\tau} \prod_{i=1}^n [N(\theta_i | \mu, \tau^2) N(y_i | \theta_i, \sigma_i^2)]$ , then:

$$\begin{aligned} p(\theta_i | \mu, \tau^2, \mathbf{y}, \boldsymbol{\sigma}^2) &\propto N(\theta_i | \mu, \tau^2) N(y_i | \theta_i, \sigma_i^2) \\ &\propto e^{-\frac{1}{2}(\frac{1}{\tau^2} + \frac{1}{\sigma_i^2})\theta_i^2 + (\frac{\mu}{\tau^2} + \frac{y_i}{\sigma_i^2})\theta_i} \\ &\Rightarrow p(\theta_i | \mu, \tau^2, \mathbf{y}, \boldsymbol{\sigma}^2) = N(\theta_i | \mu_i^*, \tau_i^{2*}) \end{aligned}$$

where,  $\tau_i^{2*} = (\frac{1}{\tau^2} + \frac{1}{\sigma_i^2})^{-1}$ , and  $\mu^* = (\frac{\mu}{\tau^2} + \frac{y_i}{\sigma_i^2})\tau_i^{2*}$

1000 samples of  $\boldsymbol{\theta}$  were drawn by using the samples of  $\mu$  and  $\tau^2$ . The mean of each  $\theta_i$  is:

$\theta_1$	$\theta_2$	$\theta_3$	$\theta_4$
1.11757159	0.81483007	-0.14736517	0.71563212
$\theta_5$	$\theta_6$	$\theta_7$	$\theta_8$
-0.02328974	0.15142220	1.71089974	0.20240492

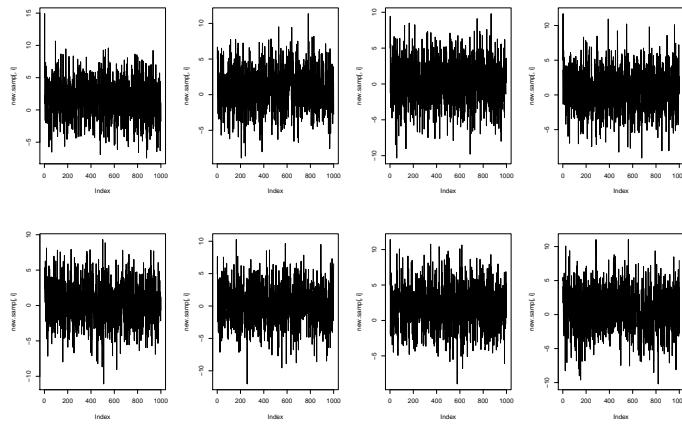
5. For Questions 5-7, consider  $\tau^2$  to be fixed and equal to  $\tau^2 = \text{median}(\tau^2)$  that you calculated in Question 3. We now use a Gibbs sampler to draw samples from the posterior distribution  $p(\mu, \boldsymbol{\theta} | \mathbf{y}, \boldsymbol{\sigma}^2)$ . Remember that  $\tau^2$  is now a known quantity in this question, so we do not need to calculate the posterior distribution for it. Give the conditional distributions of the remaining unknown parameters given the other parameters.

Solution: When  $\tau^2$  is fixed, the joint posterior of  $\mu, \boldsymbol{\theta}$  is (we denote  $\tau_0^2 = \text{median}(\tau^2)$ ):

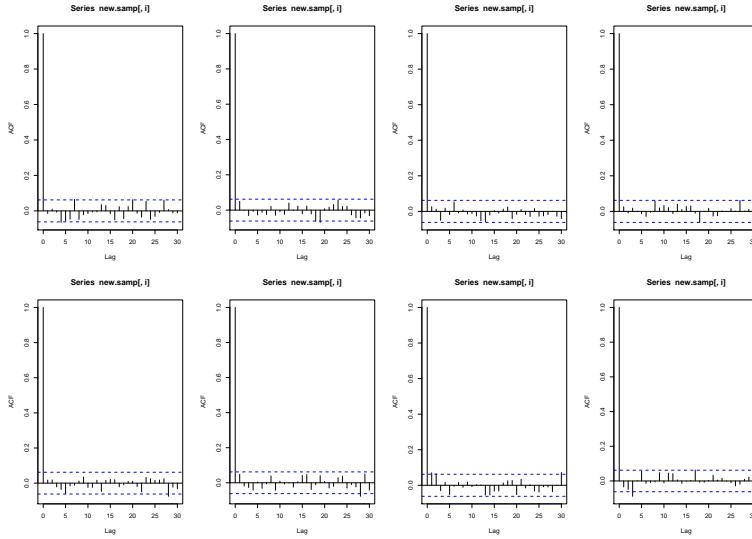
$$\begin{aligned}
p(\mu, \boldsymbol{\theta} | \mathbf{y}, \boldsymbol{\sigma}^2, \tau_0^2) &\propto p(\mu)p(\boldsymbol{\theta} | \mu, \tau_0^2)p(\mathbf{y} | \boldsymbol{\theta}, \boldsymbol{\sigma}^2) \propto \prod_{i=1}^n [N(\theta_i | \mu, \tau_0^2)p(y_i | \theta_i, \sigma_i^2)] \\
P(\theta_i | \mu, \mathbf{y}, \boldsymbol{\sigma}^2, \tau_0^2) &\propto N(\theta_i | \mu_i, \tau_0^2)p(y_i | \theta_i, \sigma_i^2) \propto e^{-\frac{1}{2}(\frac{1}{\tau_0^2} + \frac{1}{\sigma_i^2})\theta_i^2 + (\frac{\mu}{\tau_0^2} + \frac{y_i}{\sigma_i^2})\theta_i} \\
&\Rightarrow p(\theta_i | \mu, \mathbf{y}, \boldsymbol{\sigma}^2, \tau_0^2) = N(\theta_i | \xi_i, \gamma_i^2) \\
&\gamma_i^2 = (\frac{1}{\tau_0^2} + \frac{1}{\sigma_i^2})^{-1}, \quad \xi_i = (\frac{\mu}{\tau_0^2} + \frac{y_i}{\sigma_i^2})\gamma_i^2 \\
p(\mu | \boldsymbol{\theta}, \mathbf{y}, \boldsymbol{\sigma}^2, \tau_0^2) &\propto \prod_{i=1}^n N(\theta_i | \mu, \tau_0^2) \propto e^{-\frac{n}{2\tau_0^2}\mu^2 + \frac{\sum \theta_i}{\tau_0^2}\mu} \\
&\Rightarrow p(\mu | \boldsymbol{\theta}, \mathbf{y}, \boldsymbol{\sigma}^2, \tau_0^2) = N(\mu | \frac{\sum \theta_i}{n}, \frac{\tau_0^2}{n})
\end{aligned}$$

6. Use a Gibbs sampler based on the conditional distributions from the previous question to obtain 1000 samples from  $p(\mu, \boldsymbol{\theta} | \mathbf{y}, \boldsymbol{\sigma}^2)$ . Make sure that you only use samples after convergence and that your 1000 samples have low autocorrelation. Calculate the mean and posterior interval for each  $\theta_i$ .

Solution: We can plot the samples we simulated:



The ACF graphics are:



We can see from the plot of the samples, and the ACF plot that the samples have low autocorrelation and good convergence, so we use the samples we simulated to calculate the mean and posterior interval for each  $\theta_i$ .

The mean and posterior interval for each  $\theta_i$  is:

	mean	95% posterior interval
$\theta_1$	1.4707507	[-4.28861989, 7.4656290]
$\theta_2$	0.9016350	[-5.01630629, 6.6062699]
$\theta_3$	0.0796483	[-6.00374072, 6.2903041]
$\theta_4$	0.6655649	[-5.30738770, 6.4243711]
$\theta_5$	0.2933542	[-5.52278423, 6.0655314]
$\theta_6$	0.2238173	[-5.59824110, 6.1768044]
$\theta_7$	1.8908095	[-3.94686628, 7.9728499]
$\theta_8$	0.6777565	[-5.48349935, 6.9848530]

7. Use your samples of  $\boldsymbol{\theta}$  to calculate an estimate of the posterior probability that school A offers the best program.

Solution: We draw 1000 samples for each school, and repeat this procedure 1000 times. And then we got the result:

mean	median	95% confidence interval
0.159995	0.160000	[0.143000, 0.177025]

Questions 8-11 are based on the bicycle data in Table 2 below. We will focus only on the first two rows of the table (the residential streets with bike routes). We want to model the total amount of traffic  $y_i$  on each street (eg.  $y_1 = 74$ ) as follows:

$$\begin{aligned} y_i &\sim \text{Poisson}(\theta_i) \\ \theta_i &\sim \text{Gamma}(\alpha, \beta) \end{aligned}$$

8. Write out the posterior distribution of our unknown variables  $p(\alpha, \beta, \boldsymbol{\theta} | \mathbf{y})$  using a flat prior distribution for  $\alpha$  and  $\beta$ ,  $p(\alpha, \beta) \propto 1$ .

Solution:  $y_i \sim \text{Poisson}(\theta_i)$ ,  $\theta_i \sim \text{Gamma}(\alpha, \beta)$

$$p(\alpha, \beta, \boldsymbol{\theta} | \mathbf{y}) \propto p(\alpha, \beta)p(\boldsymbol{\theta} | \alpha, \beta)p(\mathbf{y} | \boldsymbol{\theta}) \propto \prod_{i=1}^n [\text{Gamma}(\theta_i | \alpha, \beta)\text{Poisson}(y_i | \theta_i)]$$

9. We want to use a Gibbs sampler to get samples from  $p(\alpha, \beta, \boldsymbol{\theta} | \mathbf{y})$ . What is the conditional distribution of  $\theta_i$  given all the other parameters?

Solution:

$$\begin{aligned}
p(\theta_i | \boldsymbol{\theta}_{-i}, \alpha, \beta, \mathbf{y}) &= p(\theta_i | \alpha, \beta, \mathbf{y}) \\
&= \text{Gamma}(\theta_i | \alpha, \beta) \text{Poisson}(y_i | \theta_i) \\
&\propto \theta_i^{\alpha-1} e^{-\beta\theta_i} \theta_i^{y_i} e^{-\theta_i} \\
&\propto \theta_i^{\alpha+y_i-1} e^{-(\beta+1)\theta_i} \\
&= \text{Gamma}(\alpha + y_i, \beta + 1)
\end{aligned}$$

10. The conditional distributions  $p(\alpha, \beta | \mathbf{y})$  and  $p(\beta | \alpha, \mathbf{y})$  are not easy to sample from, so you instead need to use a Metropolis (or Metropolis-Hastings) step to obtain samples from  $p(\alpha, \beta | \boldsymbol{\theta}, \mathbf{y})$ . Choose a proposal distribution for  $\alpha$  and a proposal distribution for  $\beta$ .

Solution:  $\alpha, \beta$  are the parameters of Gamma distribution, therefore  $\alpha$  and  $\beta$  have to be both positive. Then choosing a log-normal distribution to be the proposal distribution seems to be reasonable. i.e. We chose log-normal random walks and firstly, we will look at the independent log-normal random walks.

$$\begin{pmatrix} \log \alpha^* \\ \log \beta^* \end{pmatrix} \sim N \left( \begin{pmatrix} \log \alpha^{(l)} \\ \log \beta^{(l)} \end{pmatrix}, \begin{pmatrix} S_\alpha & 0 \\ 0 & S_\beta \end{pmatrix} \right)$$

The proposal is:

$$\begin{aligned}
q(\alpha^* | \alpha^{(l)}) &\propto \frac{1}{\alpha^*} e^{-\frac{(\ln \alpha^* - \ln \alpha^{(l)})^2}{2S_\alpha^2}}, \quad q(\alpha^{(l)} | \alpha^*) \propto \frac{1}{\alpha^{(l)}} e^{-\frac{(\ln \alpha^{(l)} - \ln \alpha^*)^2}{2S_\alpha^2}} \\
q(\beta^* | \beta^{(l)}) &\propto \frac{1}{\beta^*} e^{-\frac{(\ln \beta^* - \ln \beta^{(l)})^2}{2S_\beta^2}}, \quad q(\beta^{(l)} | \beta^*) \propto \frac{1}{\beta^{(l)}} e^{-\frac{(\ln \beta^{(l)} - \ln \beta^*)^2}{2S_\beta^2}}
\end{aligned}$$

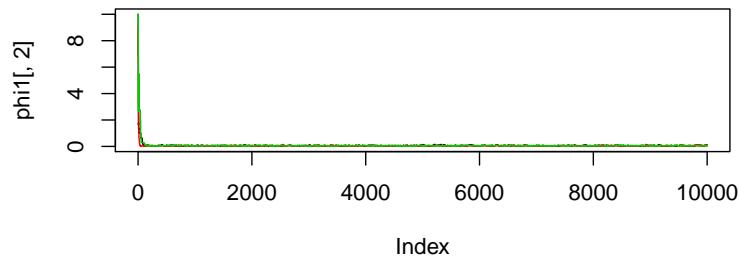
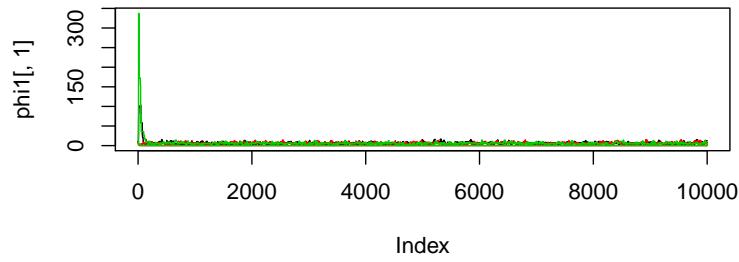
11. Combine the algorithm from Question 10 with the conditional distribution from Question 9 to form a Gibbs sampler, which iterates between sampling:
1. sampling  $\theta_i | \alpha, \beta, \mathbf{y}$  for each  $i$ .
  2. sampling  $\alpha, \beta | \boldsymbol{\theta}, \mathbf{y}$ .

Use this algorithm to obtain 1000 samples from  $p(\alpha, \beta, \boldsymbol{\theta} | \mathbf{y})$ . Make sure that you only use samples after convergence and that your 1000 samples have low autocorrelation. Calculate means and 95% posterior intervals for  $\alpha$  and  $\beta$ .

Solution: Three different standard deviation matrix were used:

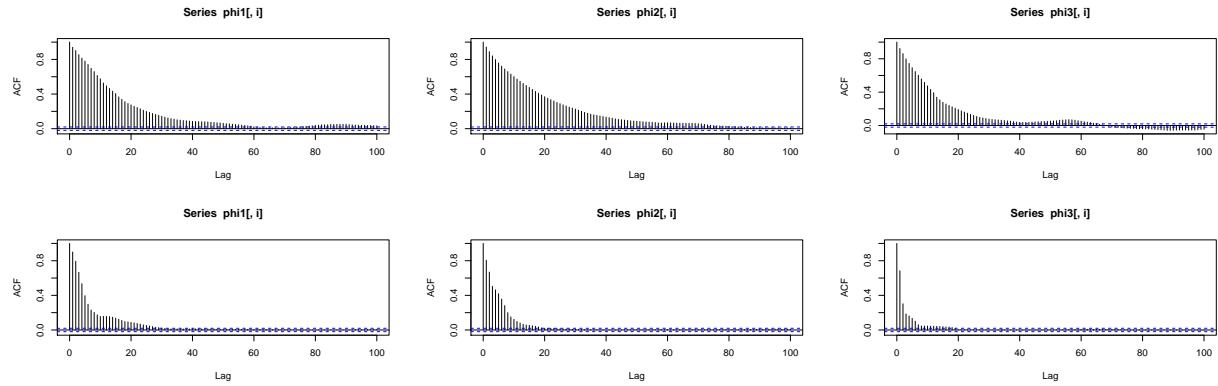
$$\begin{pmatrix} 0.35 & 0 \\ 0 & 0.35 \end{pmatrix}, \begin{pmatrix} 0.5 & 0 \\ 0 & 0.5 \end{pmatrix}, \begin{pmatrix} 0.75 & 0 \\ 0 & 0.75 \end{pmatrix}$$

We got three different chains, and here is the comparison of these three different chains:



We can see from the comparison of these three different chains that the simulation result is quite satisfactory.

Besides, we can see the different acf graphics for different chains:

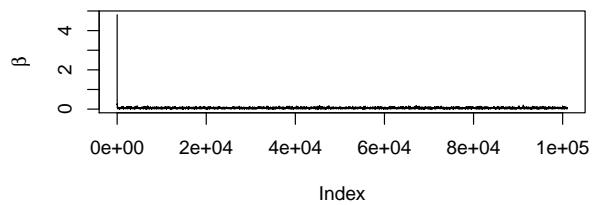
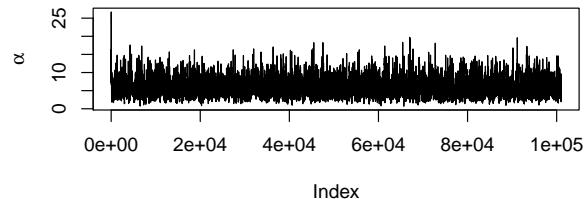


We can see the summary for these different chains:

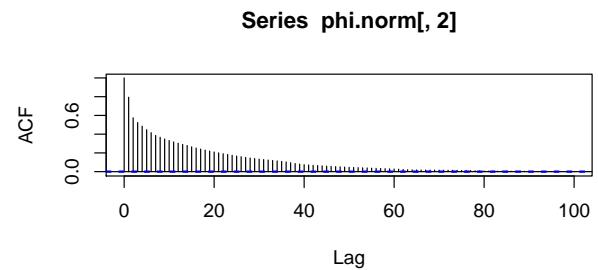
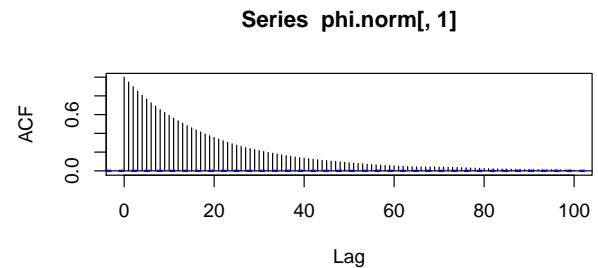
$\alpha$	median	mean	95% confidence interval
chain1	6.078	6.526	[2.516435, 3.454901]
chain2	5.574	5.971	[2.417659, 3.306146]
chian3	5.809	6.151	[2.337488, 3.213865]

$\beta$	median	mean	95% confidence interval
chain1	0.053140	0.064920	[0.02064850, 0.02972844]
chain2	0.049880	0.057450	[0.02038208, 0.02810951]
chain3	0.051240	0.056390	[0.01959811, 0.02804610]

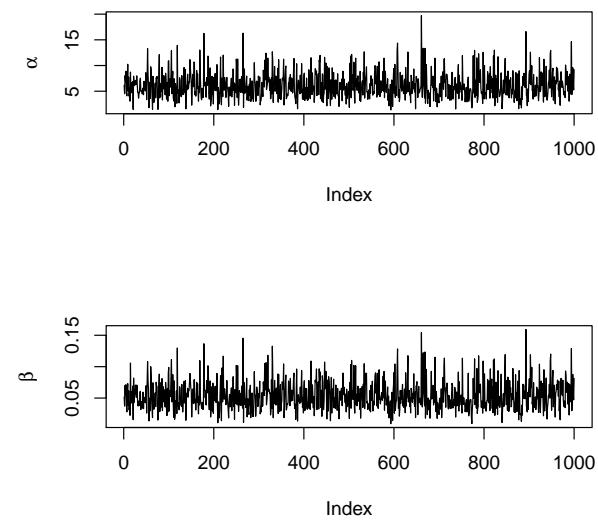
From the result before, we can see that all of the three chains are really similar and the results are all good, then a longer chain was runed using the third chain.

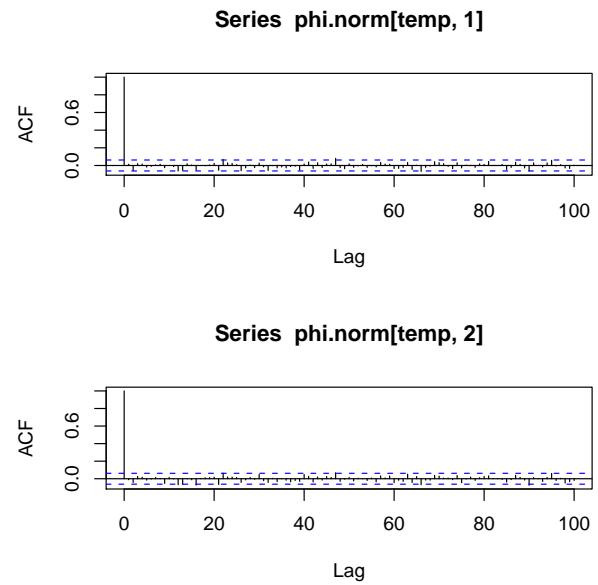


The acf is:

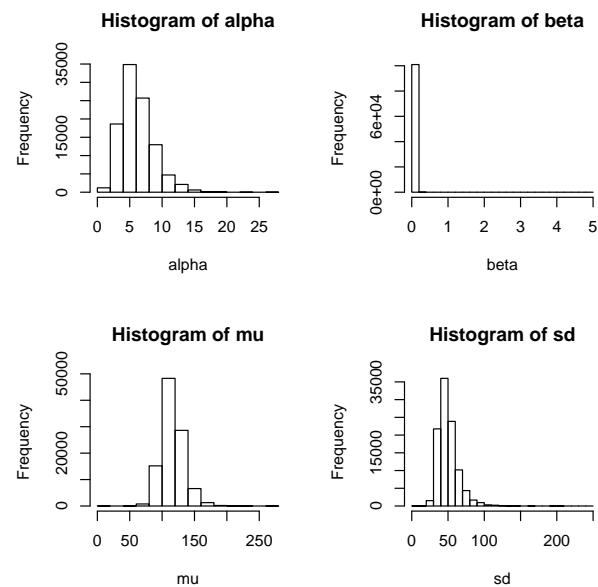


We can see that the lag is less than 100, and then we can get 1000 samples for both  $\alpha$  and  $\beta$ .

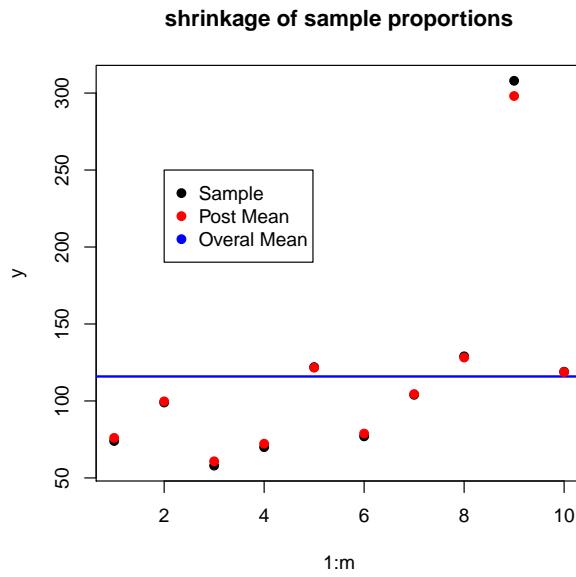




Besides, we can also examine the histogram of  $\alpha$ ,  $\beta$ , simulated mean and standard deviation.



The shrinkage situation is:

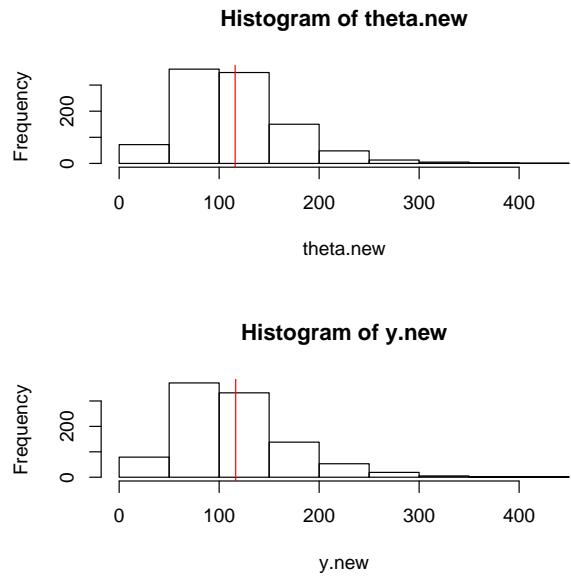


We can see that the samples we got for  $\alpha$  and  $\beta$  are convergent and approximately independent, then we can use these samples to do prediction, and here is the posterior mean and 95% confidence interval for our independent samples:

parameter	mean	95% confidence interval
$\alpha$	6.133638	[2.156641, 12.114937]
$\beta$	0.05377967	[0.0179822, 0.02791607]

12. Give a 95% predictive interval for the total traffic on a new residential street with a bike lane.

Solution: We have already got 1000 independent samples for  $\alpha$  and  $\beta$ , we can use them to predive the total traffic on a new residential street with a bike lane. The histogram for the preditive  $y$  and  $\theta$  is:



The 95% predictive interval for  $y$  is: [34.0, 253.1].

School	Estimated treatment effect, $y_j$	Standard error of estimate, $\sigma_j$
A	28	15
B	8	10
C	-3	16
D	7	11
E	-1	9
F	1	11
G	18	10
H	12	18

Table 1: *Observed effects of special preparation on SAT-V scores in eight randomized experiments. Estimates are based on separate analyses for the eight experiments. From Rubin (1981).*

Type of street	Bike route?	Counts of bicycles/other vehicles
Residential	yes	16/58, 9/90, 10/48, 13/57, 19/103, 20/57, 18/86, 17/112, 35/273, 55/64
Residential	no	12/113, 1/18, 2/14, 4/44, 9/208, 7/67, 9/29, 8/154
Fairly busy	yes	8/29, 35/415, 31/425, 19/42, 38/180, 47/675, 44/620, 44/437, 29/47, 18/462
Fairly busy	no	10/557, 43/1258, 5/499, 14/601, 58/1163, 15/700, 0/90, 47/1093, 51/1459, 32/1086
Busy	yes	60/1545, 51/1499, 58/1598, 59/503, 53/407, 68/1494, 68/1558, 60/1706, 71/476, 63/752
Busy	no	8/1248, 9/1246, 6/1596, 9/1765, 19/1290, 61/2498, 31/2346, 75/3101, 14/1918, 25/2318

Table 2: *Counts of bicycles and other vehicles in one hour in each of 10 city blocks in each of six categories. (The data for two of the residential blocks were lost.) For example, the first block had 16 bicycles and 58 other vehicles, the second had 9 bicycles and 90 other vehicles, and so on. Streets were classified as 'residential,' 'fairly busy,' or 'busy' before the data were gathered.*