



A Guide to Bayesian Modeling

Radboudumc

Author: Belias Michael

Student in MSc in Biostatistics

Athens, 20 April 2016

Bayesian Analysis

Michael Belias

April 20, 2017

Contents

1	Introduction	3
2	Common probability distributions	4
2.1	Discrete distributions	4
3	Bayesian Theory	11
3.1	History	11
3.2	Bayes theorem	12
4	Monte Carlo estimation	13
4.1	Simulation and CLT	13
4.2	Calculating probabilities	14
4.3	Monte Carlo error	14
4.4	Marginalization	15
5	Markov chains	16
5.1	Examples of Markov chains	16
5.2	Monte Carlo Example	25
6	Metropolis-Hastings	26
6.1	Proposal distribution	27
6.2	Acceptance rate	27
7	Popular Models	29
7.1	Linear Regression	30
7.2	(M)ANOVA	31
7.3	Poisson regression	32
8	Multi-level modeling	33
8.1	Hierarchical models	34
8.2	Meta analysis	52
	Bibliography	53

1 Introduction

According to the Oxford dictionary, statistics is a branch of mathematics dealing with the collection, analysis, interpretation, presentation, and organization of data. Data may be of any applied science field such as medical, finance, social, physics etc and they can be separated into 2 types quantitative and qualitative.

The typical steps of a statistical analysis are seven in general :

- 1) Define the problem
- 2) Data collection and manipulation
- 3) Explore the data
- 4) Using the above three decide the model that will be used
- 5) Fit the model
- 6) Check the model and develop if necessary
- 7) Make the final model and infer

The above steps are not distinct and in some cases there are overlaps and more steps nested. The same principles can be applied in the Bayesian Framework too.

In this tutorial we will learn:

- The bayesian intuition
- Fit the bayesian methods in simple popular statistical approaches such as:
 - (M)ANOVA
 - Linear Regression
 - Poisson regression
- Multi-level modeling
 - Hierarchical models
 - Meta analysis

2 Common probability distributions

2.1 Discrete distributions

2.1.1 Uniform

The uniform distribution is the simplest discrete probability distribution. It assigns equal probability to N different outcomes, usually represented with numbers $1, 2, \dots, N$.

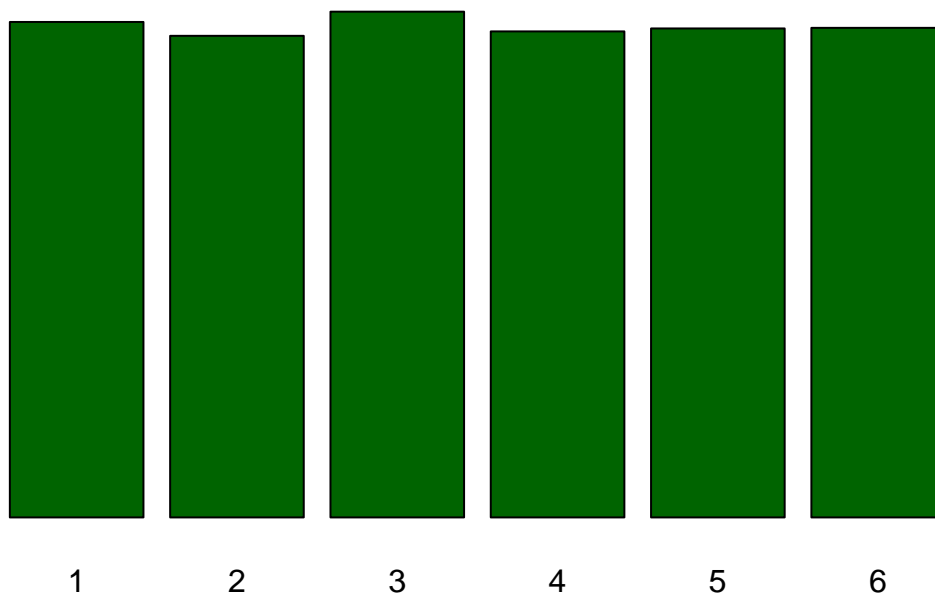
$$X \sim \text{Uniform}(N)$$

$$P(X = x|N) = 1/N \text{ for } x = 1, 2, \dots, N$$

$$E[X] = \frac{N+1}{2}$$

$$\text{Var}[X] = \frac{N^2+1}{12}$$

One common example is the outcome of throwing a fair six-sided die where $N=6$.





2.1.2 Bernoulli

The Bernoulli distribution is used for binary outcomes, such as 0 and 1. It has one parameter p , which is the probability of “success” frequently getting 1 (or any value we set). $X \sim \text{Bern}(p)$

$$P(X = x|p) = p^x(1 - p)^{1-x} \text{ for } x = 0, 1$$

$$E[X] = p$$

$$\text{Var}[X] = p(1-p)$$

One common example is the outcome of flipping a fair coin ($p = 0.5$)

2.1.3 Binomial

The binomial distribution counts the number of “successes” in n independent Bernoulli trials (each with the same probability of success). Thus if X_1, X_2, \dots, X_n are independent Bernoulli(p) random variables, then $Y = \sum_{i=1}^n X_i$ is binomial distributed.

$$Y \sim \text{Binom}(n, p)$$

$$P(Y=y|n,p) = \binom{n}{y} p^y (1-p)^{(n-y)}, \text{ for } y = 0, 1, \dots, n$$

$$E[Y] = np$$

$$\text{Var}[Y] = np(1-p)$$

$$\text{where } \binom{n}{y} = \frac{n!}{y!(n-y)!}.$$



2.1.4 Poisson

The Poisson distribution is used for counts, and arises in a variety of situations. The parameter $\lambda > 0$ is the rate at which we expect to observe the thing we are counting.

$$X \sim \text{Pois}(\lambda)$$

$$P(X = x | \lambda) = \frac{\lambda^x \exp(-\lambda)}{x!}$$

$$E[X] = \lambda$$

$$\text{Var}[X] = \lambda$$

A Poisson process is a process wherein events occur on average at rate λ , events occur one at a time, and events occur independently of each other.

Example:

Significant earthquakes occur in the Western United States approximately following a Poisson process with rate of two earthquakes per week. What is the probability there will be at least 3 earthquakes in the next two weeks? Answer: the rate per two weeks is $2 \times 2 = 4$, so let $X \sim \text{Pois}(4)$ and we want to know $P(X \geq 3) = 1 - (X \leq 2) = 1 - P(X = 0) - P(X = 1) - P(X = 2) = 1 - e^{-4} - 4e^{-4} - \frac{4^2 e^{-4}}{2} = 1 - 13e^{-4} = 0.762$

Note that $0! = 1$ by definition.



2.1.5 Geometric

The geometric distribution is the number of failures before obtaining the first success, i.e., the number of Bernoulli failures until a success is observed, such as the first head when flipping a coin. It takes values on the positive integers starting with 0 (alternatively, we could count total trials until first success, in which case we would start with 1).

$X \sim \text{Geo}(p)$

$$P(X = x|p) = p(1 - p)^x, \text{ for } x=1,2,\dots$$

$$E[X] = \frac{1-p}{p}$$

If the probability of getting a success is p , then the expected number of trials until the first success is $1/p$ and the expected number of failures until the first success is $(1 - p)/p$.



2.1.6 Negative Binomial

The negative binomial distribution extends the geometric distribution to model the number of failures before achieving the r th success. It takes values on the positive integers starting with 0.

$$Y \sim \text{NegBinom}(r, p)$$

$$P(Y = y | n, p) = \binom{r+y-1}{y} p^r (1-p)^y \text{ for } y=1, 2, \dots$$

$$E[Y] = \frac{r(1-p)}{p}$$

$$\text{Var}[Y] = \frac{r(1-p)}{p^2}$$

Note that the geometric distribution is a special case of the negative binomial distribution where $r = 1$. Because $0 < p < 1$, we have $E[Y] < \text{Var}[Y]$. This makes the negative binomial a popular alternative to the Poisson when modeling counts with high variance (recall, that the mean equals the variance for Poisson distributed variables).



2.1.7 Multinomial

Another generalization of the Bernoulli and the binomial is the multinomial distribution, which is like a binomial when there are more than two possible outcomes. Suppose we have n trials and there are k different possible outcomes which occur with probabilities p_1, p_2, \dots, p_k . For example, we are rolling a six-sided die that might be loaded so that the sides are not equally likely, then n is the total number of rolls, $k = 6$, p_1 is the probability of rolling a one, and we denote by x_1, x_2, \dots, x_6 a possible outcome for the number of times we observe rolls of each of one through six, where $\sum_{i=1}^6 x_i = n$ and $\sum_{i=1}^6 p_i = 1$

3 Bayesian Theory

3.1 History

Bayesian statistics are based on the homonymous Bayes' theorem or rule, invented by Thomas Bayes, which was a british reverend the 1740s . His primary field of studying was theology but Bayes was also “amateur” mathematician. He was influenced by David Hume a philosopher teacher while his studies in Edinburgh proposing that we can only rely on what we learn from experience. The probabilities as a mathematical field these days were just emerging being able to solve simple problems like *what is the probability of observing an effect given a cause?* but not the inverse $P(\text{cause} | \text{effect})$. Bayes gave a simple example of tossing balls on a table and recording where they stop (to the left or to the right side of the table), noting that the more balls are thrown, the better we may infer if the ball-tossing is biased to a side. This is nowadays called a learning process and although it was a remarkable finding Bayes forgot it in a drawer (!) until his death. Richard Price found it and after studying his papers for 2 years and making some corrections he finally published **An Essay toward solving a Problem in the Doctrine of Chances**. 1763.

Still the theorem was just an example not having the final form of today and even after this publication no-one really continued the development except of Laplace, who was trying to solve an astronomical problem , studied Price's paper developed a first version of what we now call Bayes theorem. The reception of Laplace's proposal was slightly hostile due to the inherent challenges such as the equal prior probabilities, being subjective and the serious technical computational problems in practice, which is still a great issue .



3.2 Bayes theorem

Bayes theorem is calculating the probability event given prior knowledge of conditions that might be related to the event. Bayes' theorem is stated mathematically as the following equation (Efron, 2013) :

$$P(A | B) = \frac{P(B|A)P(A)}{P(B)}$$

$$P(A | B) = \frac{P(B|A)P(A)}{P(B)}$$

where A and B are events and $P(B) \neq 0$.

- $P(A)$ and $P(B)$ are the probabilities of observing A and B without regard to each other.
- $P(A | B)$, a conditional probability, is the probability of observing event A given that B is true.
- $P(B | A)$ is the probability of observing event B given that A is true.

This is the basis of Bayesian inference which is a particular approach to statistical inference, with it's own interpretation and When applied, the probabilities involved in Bayes' theorem may have different probability interpretations. With the Bayesian probability interpretation the theorem expresses how a subjective degree of belief should rationally change to account for availability of related evidence. Bayesian inference is fundamental to Bayesian statistics.

4 Monte Carlo estimation

4.1 Simulation and CLT

Before we learn how to simulate from complicated posterior distributions, let's review some of the basics of Monte Carlo estimation. Monte Carlo estimation refers to simulating hypothetical draws from a probability distribution in order to calculate important quantities. These quantities might include the mean, the variance, the probability of some event, or quantiles of the distribution. All of these calculations involve integration, which except for the simplest distributions, can be very difficult or impossible.

Suppose we have a random variable $\hat{\theta}$ that follows a $\text{Gamma}(a, b)$. Let's say $a=2$ and $b=1/3$, where b is the rate parameter. To calculate the mean of this distribution, we would need to compute the following integral

$$E(\theta) = \int_0^{\infty} \theta f(\theta) d\theta = \int_0^{\infty} \theta \frac{b^a}{\Gamma(a)} \theta^{a-1} e^{-b\theta} d\theta$$

It is possible to compute this integral, and the answer is a/b (6 in this case). However, we could verify this answer through Monte Carlo estimation. To do so, we would simulate a large number of draws (call them θ_i^* for $i = 1, 2, \dots, m$) from this gamma distribution and calculate their sample mean. Why can we do this? Recall from the previous course that if we have a random sample from a distribution, the average of those samples converges in probability to the true mean of that distribution by the Law of Large Numbers. Furthermore, by the Central Limit Theorem (CLT), this sample mean $\bar{\theta}^* = \frac{1}{m} \sum_{i=1}^m \theta_i^*$ approximately follows a normal distribution with mean $E(\theta)$ and variance $\text{Var}(\theta)/m$. The theoretical variance of θ is the following integral:

$$\text{Var}(\theta) = \int_0^{\infty} (\theta - E(\theta))^2 f(\theta) d\theta$$

Just like we did with the mean, we could approximate this variance with the sample variance $\frac{1}{m} \sum_{i=1}^m (\theta_i^* - \bar{\theta}^*)^2$



4.2 Calculating probabilities

This method can be used to calculate many different integrals. Say $h(\theta)$ is any function and we want to calculate $\int h(\theta)p(\theta)d\theta$. This integral is precisely what is meant by $E(h(\theta))$, so we can conveniently approximate it by taking the sample mean of $h(\theta_i^*)$. That is, we apply the function h to each simulated sample from the distribution, and take the average of all the results.

One extremely useful example of an h function is the indicator $I_A(\theta)$ where A is some logical condition about the value of θ . To demonstrate, suppose $h(\theta) = I_{[\theta < 5]}(\theta)$, which will give a 1 if $\theta < 5$ and 0 otherwise.

The $E(h(\theta)) = \int_0^\infty I_{[\theta < 5]}(\theta)p(\theta)d\theta = \int_0^5 1 \cdot p(\theta)d\theta + \int_5^\infty 0 \cdot p(\theta)d\theta = P(\theta < 5)$. This means we can approximate the probability that $\theta < 5$ by drawing many samples θ_i^* , and approximating this integral with $\frac{1}{m} \sum_{i=1}^m I_{\theta^* < 5}(\theta_i^*)$. This expression is simply counting how many of those samples come out to be less than 55, and dividing by the total number of simulated samples. So simple!

Likewise, we can approximate quantiles of a distribution. If we are looking for the value ζ such that $P(\theta < \zeta) = 0.9$, we simply arrange the samples θ_i^* in ascending order and find the smallest drawn value that is greater than 90% of the others.

4.3 Monte Carlo error

How good is an approximation by Monte Carlo sampling? Again we can turn to the CLT, which tells us that the variance of our estimate is controlled in part by m . For a better estimate, we want larger m .

For example, if we seek $E(\theta)$, then the sample mean $\bar{\theta}^*$ approximately follows a normal distribution with mean $E(\theta)$ and variance $\text{Var}(\theta)/m$. The variance tells us how far our estimate might be from the true value. One way to approximate $\text{Var}(\theta)$ is to replace it with the sample variance. The standard deviation of our Monte Carlo estimate is the square root of that, or the sample standard deviation divided by \sqrt{m} .



If m is large, it is reasonable to assume that the true value will likely be within about two standard deviations of your Monte Carlo estimate.

4.4 Marginalization

We can also obtain Monte Carlo samples from hierarchical models. As a simple example, let's consider a binomial random variable where $y \mid \phi \sim \text{Bin}(10, \phi)$, and further suppose ϕ is random (as if it had a prior) and is distributed beta $\phi \sim \text{Beta}(2, 2)$. Given any hierarchical model, we can always write the joint distribution of y and ϕ as $p(y, \phi) = p(y \mid \phi)p(\phi)$ using the chain rule of probability. To simulate from this joint distribution, repeat these steps for a large number m :

- Simulate ϕ_i^* from its $\text{Beta}(2, 2)$ distribution
- Given the drawn ϕ_i^* , simulate y_i^* from $\text{Bin}(10, \phi_i^*)$

This will produce m independent pairs of (y_i^*, ϕ_i^*) drawn from their joint distribution. One major advantage of Monte Carlo simulation is that marginalizing is easy. Calculating the marginal distribution of y , $p(y) = \int_0^1 p(y, \phi) d\phi$ might be challenging. But if we have draws from the joint distribution, we can just discard the ϕ_i^* rows and use the y_i^* as samples from their marginal distribution. This is also called the prior predictive distribution introduced in the previous course.

In the next segment, we will demonstrate some of these principles. Remember, we do not yet know how to sample from the complicated posterior distributions introduced in the previous lesson. But once we learn that, we will be able to use the principles from this lesson to make approximate inferences from those posterior distributions.

5 Markov chains

Definition If we have a sequence of random variables X_1, X_2, \dots, X_n where the indices $1, 2, \dots, n$ represent successive points in time, we can use the chain rule of probability to calculate the probability of the entire sequence:

$$p(X_{t+1}|X_t, X_{t-1}, \dots, X_2, X_1) = p(X_{t+1}|X_t)$$

Markov chains simplify this expression by using the *Markov assumption*. The assumption is that given the entire past history, the probability distribution for the random variable at the next time step only depends on the current variable. Mathematically, the assumption is written like this:

$$p(X_{t+1}|X_t, X_{t-1}, \dots, X_2, X_1) = p(X_{t+1}|X_t)$$

for all $t=2, \dots, n$. Under this assumption, we can write the first expression as $p(X_1, X_2, \dots, X_n) = p(X_1) \cdot p(X_2|X_1) \cdot p(X_3|X_2) \cdot p(X_4|X_3) \cdot \dots \cdot p(X_n|X_{n-1})$,

which is much simpler than the original. It consists of an initial distribution for the first variable, $P(X_1)$, and $n - 1$ transition probabilities. We usually make one more assumption: that the transition probabilities do not change with time. Hence, the transition from time t to time $t+1$ depends only on the value of X_t .

5.1 Examples of Markov chains

5.1.1 Discrete Markov chain

Suppose you have a secret number (make it an integer) between 1 and 5. We will call it your initial number at *step 1*. Now for each time step, your secret number will change according to the following rules:

1. Flip a coin.
- 2.



- If the coin turns up heads, then increase your secret number by one (5 increases to 1).
- If the coin turns up tails, then decrease your secret number by one (1 decreases to 5).

3. Repeat n times, and record the evolving history of your secret number.

Before the experiment, we can think of the sequence of secret numbers as a sequence of random variables, each taking on a value in $\{1,2,3,4,5\}$. Assume that the coin is fair, so that with each flip, the probability of heads and tails are both 0.5.

Does this game qualify as a true Markov chain? Suppose your secret number is currently 4 and that the history of your secret numbers is $(2,1,2,3)$. What is the probability that on the next step, your secret number will be 5? What about the other four possibilities? Because of the rules of this game, the probability of the next transition will depend only on the fact that your current number is 4. The numbers further back in your history are irrelevant, so this is a Markov chain.

This is an example of a discrete Markov chain, where the possible values of the random variables come from a discrete set. Those possible values (secret numbers in this example) are called states of the chain. The states are usually numbers, as in this example, but they can represent anything. In one common example, the states describe the weather on a particular day, which could be labeled as 1-fair, 2-poor.

5.1.2 Random walk (continuous)

Now let's look at a continuous example of a Markov chain. Say $X_t=0$ and we have the following transition model: $p(X_{t+1}|X_t = x_t) = N(x_t, 1)$. That is, the probability distribution for the next state is Normal with variance 1 and mean equal to the current state. This is often referred to as a "random walk." Clearly, it is a Markov chain because the transition to the next state X_{t+1} only depends on the current state X_t .

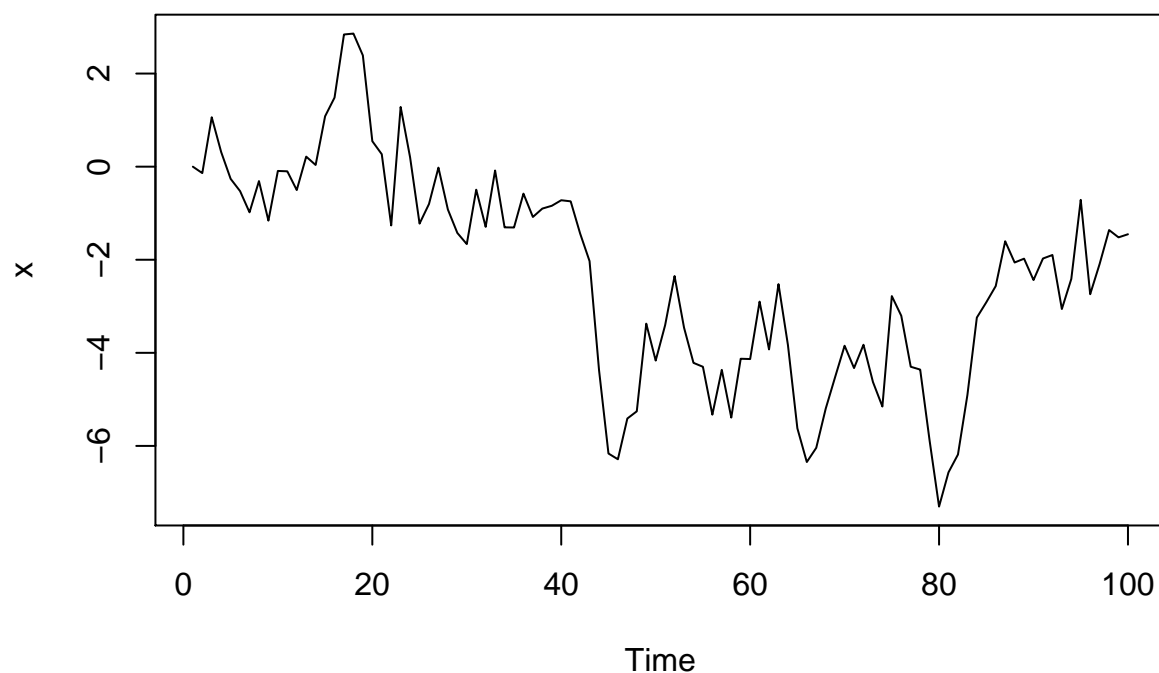
R-code

```
set.seed(34)

n = 100
x = numeric(n)

for (i in 2:n) {
  x[i] = rnorm(1, mean=x[i-1], sd=1.0)
}

plot.ts(x)
```

**5.1.3 Transition matrix**

Let's return to our example of the discrete Markov chain. If we assume that transition probabilities do not change with time, then there are a total of $5^2 = 25$ potential transition probabilities. Potential transition probabilities would be from *State 1* to *State 2*, *State 1* to *State 3*, and so forth. These transition probabilities can be arranged

into a matrix Q :

$$Q = \begin{pmatrix} 0 & .5 & 0 & 0 & .5 \\ .5 & 0 & .5 & 0 & 0 \\ 0 & .5 & 0 & .5 & 0 \\ 0 & 0 & .5 & 0 & .5 \\ .5 & 0 & 0 & .5 & 0 \end{pmatrix}$$

where the transitions from *State 1* are in the first row, the transitions from *State 2* are in the second row, etc. For example, the probability $p(X_{t+1} = 5 \mid X_t = 4)$ can be found in the fourth row, fifth column.

The transition matrix is especially useful if we want to find the probabilities associated with multiple steps of the chain. For example, we might want to know $p(X_{t+2} = 3 \mid X_t = 1)$, the probability of your secret number being 3 two steps from now, given that your number is currently 1. We can calculate this as $\sum_{k=1}^5 p(X_{t+2} = 3 \mid X_{t+1} = k) \cdot p(X_{t+1} = k \mid X_t = 1)$, which conveniently is found in the first row and third column of Q^2 .

R-code

```
Q = matrix(c(0.0, 0.5, 0.0, 0.0, 0.5,
             0.5, 0.0, 0.5, 0.0, 0.0,
             0.0, 0.5, 0.0, 0.5, 0.0,
             0.0, 0.0, 0.5, 0.0, 0.5,
             0.5, 0.0, 0.0, 0.5, 0.0),
           nrow=5, byrow=TRUE)

Q %*% Q # Matrix multiplication in R. This is Q^2.

##      [,1] [,2] [,3] [,4] [,5]
## [1,] 0.50 0.00 0.25 0.25 0.00
## [2,] 0.00 0.50 0.00 0.25 0.25
```



```
## [3,] 0.25 0.00 0.50 0.00 0.25
## [4,] 0.25 0.25 0.00 0.50 0.00
## [5,] 0.00 0.25 0.25 0.00 0.50
```

5.1.4 Stationary distribution

Suppose we want to know the probability distribution of the your secret number in the distant future, say $p(X_{t+h}|X_t)$ where h is a large number. Let's calculate this for a few different values of h .

```
Q5 = Q %*% Q %*% Q %*% Q %*% Q # h=5 steps in the future
round(Q5, 3)
```

```
##      [,1] [,2] [,3] [,4] [,5]
## [1,] 0.062 0.312 0.156 0.156 0.312
## [2,] 0.312 0.062 0.312 0.156 0.156
## [3,] 0.156 0.312 0.062 0.312 0.156
## [4,] 0.156 0.156 0.312 0.062 0.312
## [5,] 0.312 0.156 0.156 0.312 0.062
```

```
Q10 = Q %*% Q %*% Q %*% Q %*% Q %*% Q %*% Q %*% Q %*% Q # h=10 steps
round(Q10, 3)
```

```
##      [,1] [,2] [,3] [,4] [,5]
## [1,] 0.248 0.161 0.215 0.215 0.161
## [2,] 0.161 0.248 0.161 0.215 0.215
## [3,] 0.215 0.161 0.248 0.161 0.215
## [4,] 0.215 0.215 0.161 0.248 0.161
## [5,] 0.161 0.215 0.215 0.161 0.248
```

```
Q30 = Q
for (i in 2:30) {
  Q30 = Q30 %*% Q
}
```



```

}
round(Q30, 3) # h=30 steps in the future

##      [,1]  [,2]  [,3]  [,4]  [,5]
## [1,] 0.201 0.199 0.200 0.200 0.199
## [2,] 0.199 0.201 0.199 0.200 0.200
## [3,] 0.200 0.199 0.201 0.199 0.200
## [4,] 0.200 0.200 0.199 0.201 0.199
## [5,] 0.199 0.200 0.200 0.199 0.201

```

Notice that as the future horizon gets more distant, the transition distributions appear to converge. The state you are currently in becomes less important in determining the more distant future. If we let h get really large, and take it to the limit, all the rows of the long-range transition matrix will become equal to $(.2,.2,.2,.2,.2)$. That is, if you run the Markov chain for a very long time, the probability that you will end up in any particular state is $1/5=.2$ for each of the five states. These long-range probabilities are equal to what is called the stationary distribution of the Markov chain.

The stationary distribution of a chain is the initial state distribution for which performing a transition will not change the probability of ending up in any given state. That is,

```

c(0.2, 0.2, 0.2, 0.2, 0.2) %*% Q

##      [,1] [,2] [,3] [,4] [,5]
## [1,]  0.2  0.2  0.2  0.2  0.2

```

One consequence of this property is that once a chain reaches its stationary distribution, the stationary distribution will remain the distribution of the states thereafter.

We can also demonstrate the stationary distribution by simulating a long chain from this example.

```
n = 5000
```



```
x = numeric(n)
x[1] = 1 # fix the state as 1 for time 1
for (i in 2:n) {
  x[i] = sample.int(5, size=1, prob=Q[x[i-1],]) # draw the next state from
}
```

Now that we have simulated the chain, let's look at the distribution of visits to the five states.

```
table(x) / n

## x
##      1      2      3      4      5
## 0.1996 0.2020 0.1980 0.1994 0.2010
```

The overall distribution of the visits to the states is approximately equal to the stationary distribution.

As we have just seen, if you simulate a Markov chain for many iterations, the samples can be used as a Monte Carlo sample from the stationary distribution. This is exactly how we are going to use Markov chains for Bayesian inference. In order to simulate from a complicated posterior distribution, we will set up and run a Markov chain whose stationary distribution is the posterior distribution.

It is important to note that the stationary distribution doesn't always exist for any given Markov chain. The Markov chain must have certain properties, which we won't discuss here. However, the Markov chain algorithms we'll use in future lessons for Monte Carlo estimation are guaranteed to produce stationary distributions.

5.1.5 Continuous example

The continuous random walk example we gave earlier does not have a stationary distribution. However, we can modify it so that it does have a stationary distribution.

Let the transition distribution be $p(X_{t+1}|X_t = x_t) = N(\phi x_t, 1)$ where $-1 < \phi < 1$. That is, the probability distribution for the next state is Normal with variance 1 and mean equal to ϕ times the current state. As long as ϕ is between -1 and 1 , then the stationary distribution will exist for this model.

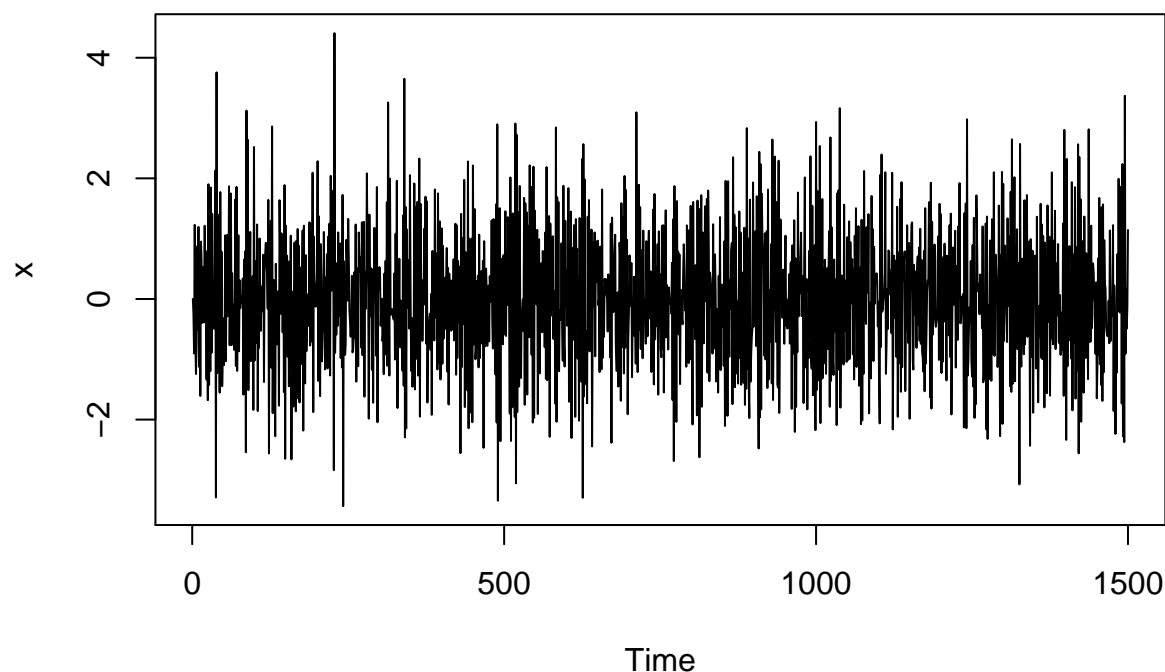
Let's simulate this chain for $\phi = -0.6$.

```
set.seed(38)

n = 1500
x = numeric(n)
phi = -0.6

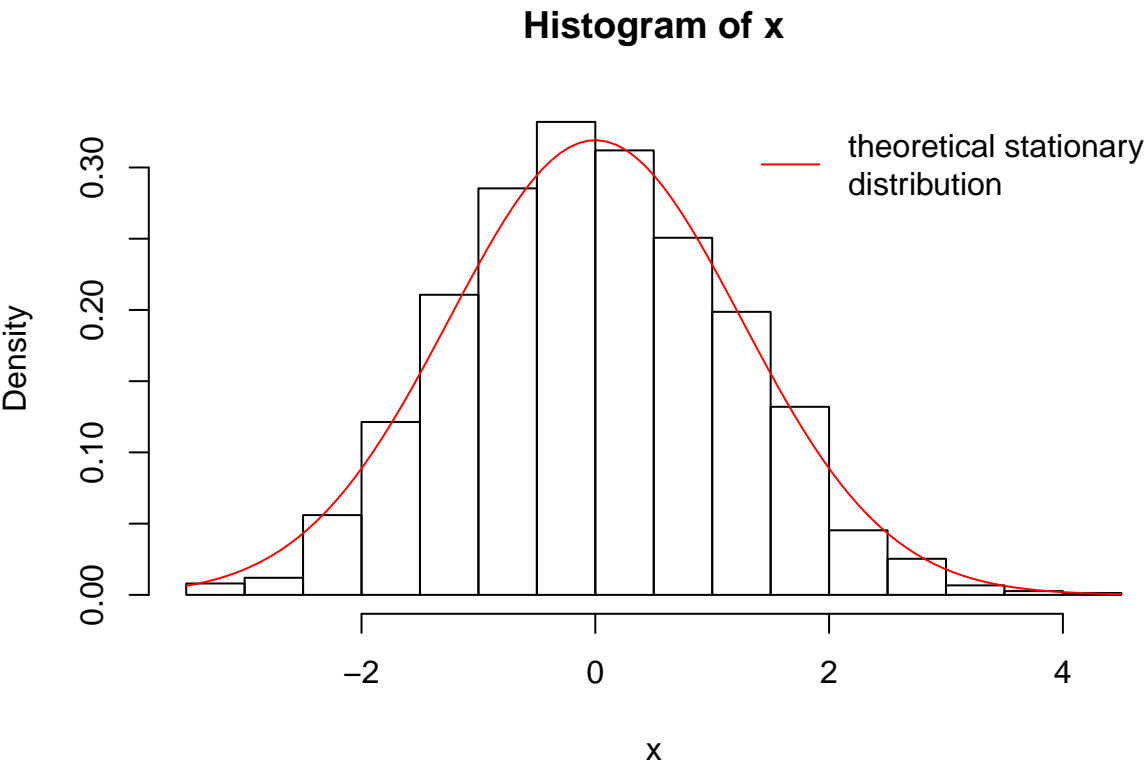
for (i in 2:n) {
  x[i] = rnorm(1, mean=phi*x[i-1], sd=1.0)
}

plot.ts(x)
```



The theoretical stationary distribution for this chain is normal with mean 0 and vari-

ance $1/(1-\phi^2)$ which in our example approximately equals 1.562 . Let's look at a histogram of our chain and compare that with the theoretical stationary distribution.



It appears that the chain has reached the stationary distribution. Therefore, we could treat this simulation from the chain like a Monte Carlo sample from the stationary distribution, a normal with mean 0 and variance 1.562 .

Because most posterior distributions we will look at are continuous, our Monte Carlo simulations with Markov chains will be similar to this example.



5.2 Monte Carlo Example

6 Metropolis-Hastings

Metropolis-Hastings is an algorithm that allows us to sample from a generic probability distribution (which we will call the target distribution), even if we do not know the normalizing constant. To do this, we construct and sample from a Markov chain whose stationary distribution is the target distribution. It consists of picking an arbitrary starting value, and iteratively accepting or rejecting candidate samples drawn from another distribution, one that is easy to sample.

Let's say we wish to produce samples from a target distribution $p(\theta) \propto g(\theta)$ where we don't know the normalizing constant (since $\int g(\theta)d\theta$ is hard or impossible to compute), so we only have $g(\theta)$ to work with. The Metropolis-Hastings algorithm proceeds as follows.

1. Select an initial value θ_0 .
2. For $i=1, \dots, m$ repeat the following steps:
 - Draw a candidate sample θ^* from a proposal distribution $q(\theta^* | \theta_{i-1})$ (more on this later). *Compute the ratio $\alpha = \frac{g(\theta^*)/q(\theta^*|\theta_{i-1})}{g(\theta_{i-1})/q(\theta_{i-1}|\theta^*)} = \frac{g(\theta^*)q(\theta_{i-1}|\theta^*)}{g(\theta_{i-1})q(\theta^*|\theta_{i-1})}$.

If $\alpha \geq 1$, then set $\theta_i = \theta^*$. If $\alpha < 1$, then set $\theta_i = \theta^*$ with probability α , or $\theta_i = \theta_{i-1}$ with probability $1-\alpha$.

Steps 2b and 2c act as a correction since the proposal distribution is not the target distribution. At each step in the chain, we draw a candidate and decide whether to “move” the chain there or remain where we are. If the proposed move to the candidate is “advantageous,” ($\alpha \geq 1$) we “move” there and if it is not “advantageous,” we still might move there, but only with probability α . Since our decision to “move” to the candidate only depends on where the chain currently is, this is a Markov chain.



6.1 Proposal distribution

One careful choice we must make is the candidate generating distribution $q(\theta^* \mid \theta_{i-1})$. It may or may not depend on the previous iteration's value of θ . One example where it doesn't depend on the previous value would be if $q(\theta)$ is always the same distribution. If we use this option, $q(\theta)$ should be as similar as possible to $p(\theta)$.

Another popular option, one that does depend on the previous iteration, is Random-Walk Metropolis-Hastings. Here, the proposal distribution is centered on θ_{i-1} . For instance, it might be a normal distribution with mean θ_{i-1} . Because the normal distribution is symmetric, this example comes with another advantage: $q(\theta^* \mid \theta_{i-1}) = q(\theta_{i-1} \mid \theta^*)$, causing it to cancel out when we calculate α . Thus, in Random-Walk Metropolis-Hastings where the candidate is drawn from a normal with mean θ_{i-1} and constant variance, the acceptance ratio is $\alpha = \frac{g(\theta^*)}{g(\theta_{i-1})}$.

6.2 Acceptance rate

Clearly, not all candidate draws are accepted, so our Markov chain sometimes “stays” where it is, possibly for many iterations. How often you want the chain to accept candidates depends on the type of algorithm you use. If you approximate $p(\theta)$ with $q(\theta^*)$ and always draw candidates from that, accepting candidates often is good; it means $q(\theta^*)$ is approximating $p(\theta)$ well. However, you still may want q to have a larger variance than p and see some rejection of candidates as an assurance that q is covering the space well.

As we will see in coming examples, a high acceptance rate for the Random-Walk Metropolis-Hastings sampler is not a good thing. If the random walk is taking too small of steps, it will accept often, but will take a very long time to fully explore the posterior. If the random walk is taking too large of steps, many of its proposals will have low probability and the acceptance rate will be low, wasting many draws. Ideally, a random walk sampler should accept somewhere between 23% and 50% of



the candidates proposed.

In the next segment, we will see a demonstration of this algorithm used in a discrete case, where we can show mathematically that the Markov chain converges to the target distribution. In the following segment, we will demonstrate coding a Random-Walk Metropolis-Hastings algorithm in R to solve one of the problems from the end of Lesson 2.

7 Popular Models



7.1 Linear Regression



7.2 (M)ANOVA



7.3 Poisson regression

8 Multi-level modeling



8.1 Hierarchical models

8.1.1 Data

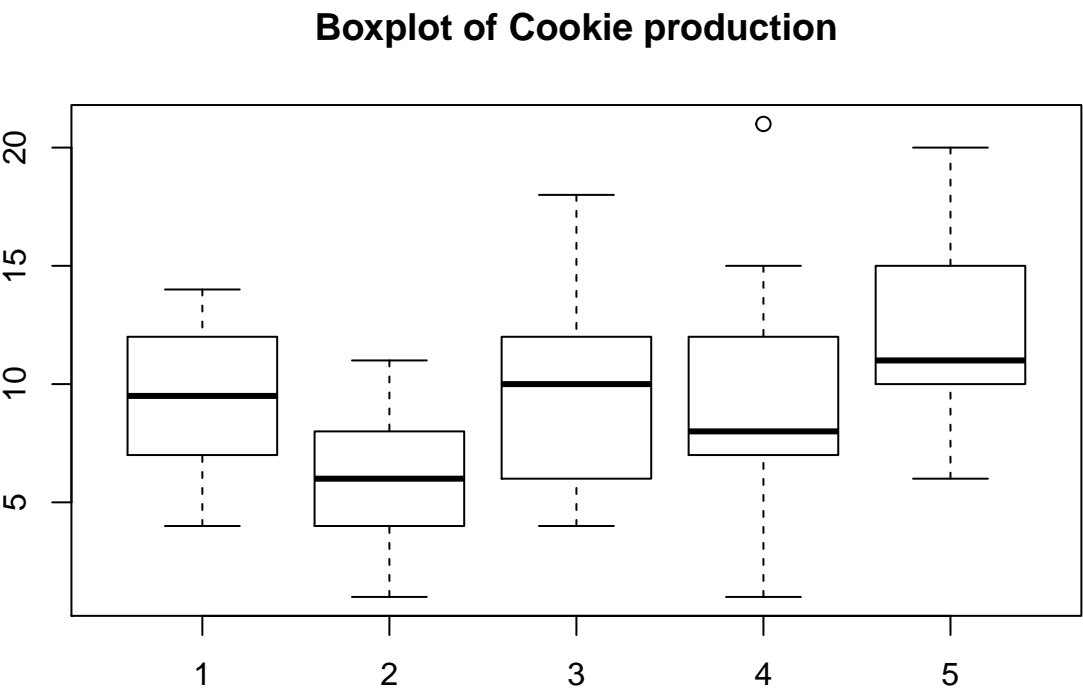
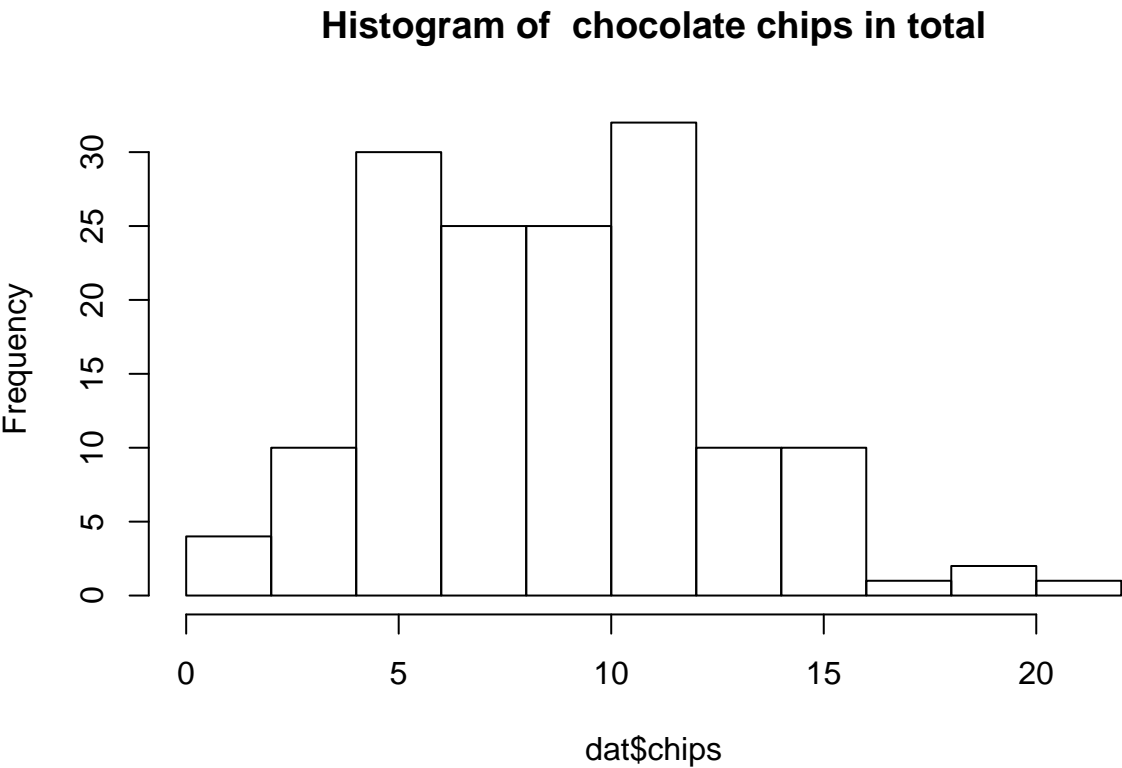
Let’s fit our hierarhical model for counts of chocolate chips. The data can be found in

Table 1: First 10 values

chips	location
12	1
12	1
6	1
13	1
12	1
12	1
9	1
10	1
7	1
14	1

```
##
## 1 2 3 4 5
## 30 30 30 30 30
```

We can also visualize the distribution of chips by location.



8.1.2 Prior predictive checks

Before implementing the model, we need to select prior distributions for α and β , the hyperparameters governing the gamma distribution for the λ parameters. First, think



about what the λ 's represent. For location j , λ_j is the expected number of chocolate chips per cookie. Hence, α and β control the distribution of these means between locations. The mean of this gamma distribution will represent the overall mean of number of chips for all cookies. The variance of this gamma distribution controls the variability between locations. If this is high, the mean number of chips will vary widely from location to location. If it is small, the mean number of chips will be nearly the same from location to location.

To see the effects of different priors on the distribution of λ 's, we can simulate. Suppose we try independent exponential priors for α and β .

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	0.021	2.983	9.852	61.127	29.980	4858.786
##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	0.1834	3.3663	8.5488	41.8137	22.2219	2865.6461

After simulating from the priors for α and β , we can use those samples to simulate further down the hierarchy:

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	0.000	1.171	7.667	83.062	28.621	11005.331

Or for a prior predictive reconstruction of the original data set:

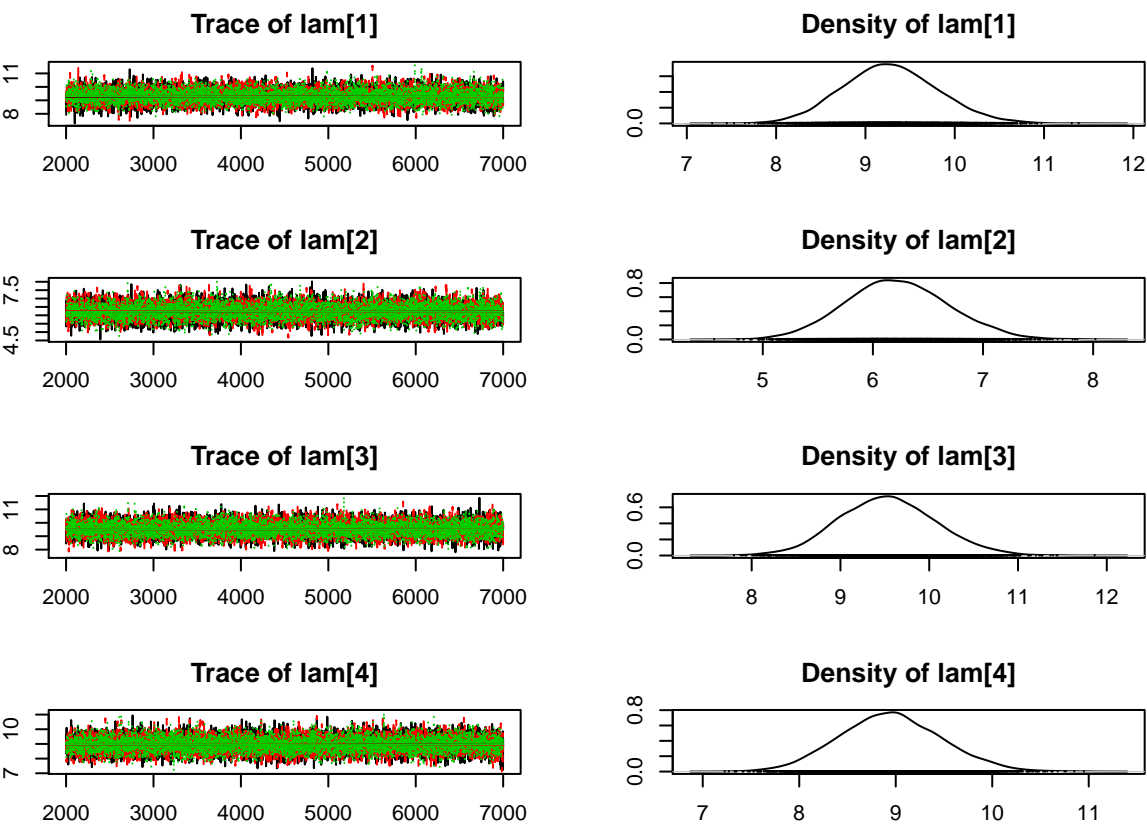
[1]	66.444084	9.946688	6.028319	15.922568	47.978587	
[1]	63	58	64	63	70	62 61 48 71 73 70 77 66 60 72 77 69 62 66 71 49 80 66
[24]	75	74	55	62	90	65 57 12 9 7 10 12 10 11 7 14 13 9 6 6 13 7 10
[47]	12	9	9	10	7	8 6 9 7 10 13 13 8 12 6 10 3 6 7 4 6 7 5
[70]	5	4	3	6	2	8 4 8 4 5 7 1 4 5 3 8 8 3 1 7 3 16 14
[93]	13	17	17	12	13	13 16 16 15 14 11 10 13 17 16 19 16 17 15 16 7 17 21
[116]	16	12	15	14	13	52 44 51 46 39 40 40 44 46 59 45 49 58 42 31 52 43 47
[139]	53	41	48	57	35	60 51 58 36 34 41 59

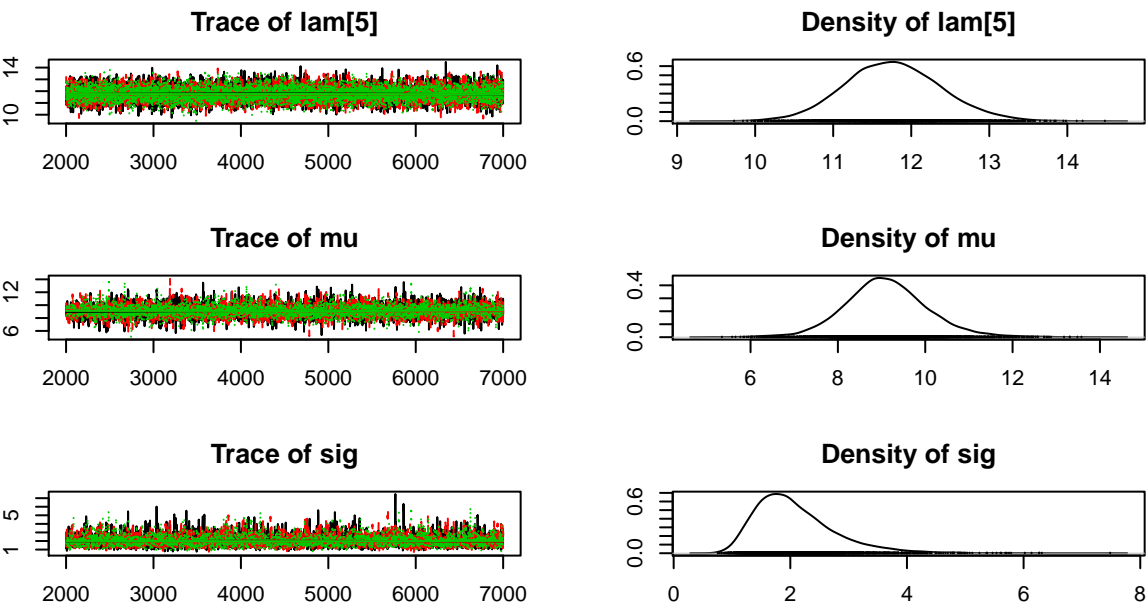
Because these priors have high variance and are somewhat noninformative, they pro-

duce unrealistic predictive distributions. Still, enough data would overwhelm the prior, resulting in useful posterior distributions. Alternatively, we could tweak and simulate from these prior distributions until they adequately represent our prior beliefs. Yet another approach would be to re-parameterize the gamma prior, which we'll demonstrate as we fit the model.

```
Compiling model graph
  Resolving undeclared variables
  Allocating nodes
Graph information:
  Observed stochastic nodes: 150
  Unobserved stochastic nodes: 7
  Total graph size: 319
```

Initializing model





Potential scale reduction factors:

	Point est.	Upper C.I.
lam[1]	1	1.00
lam[2]	1	1.00
lam[3]	1	1.00
lam[4]	1	1.00
lam[5]	1	1.00
mu	1	1.00
sig	1	1.01

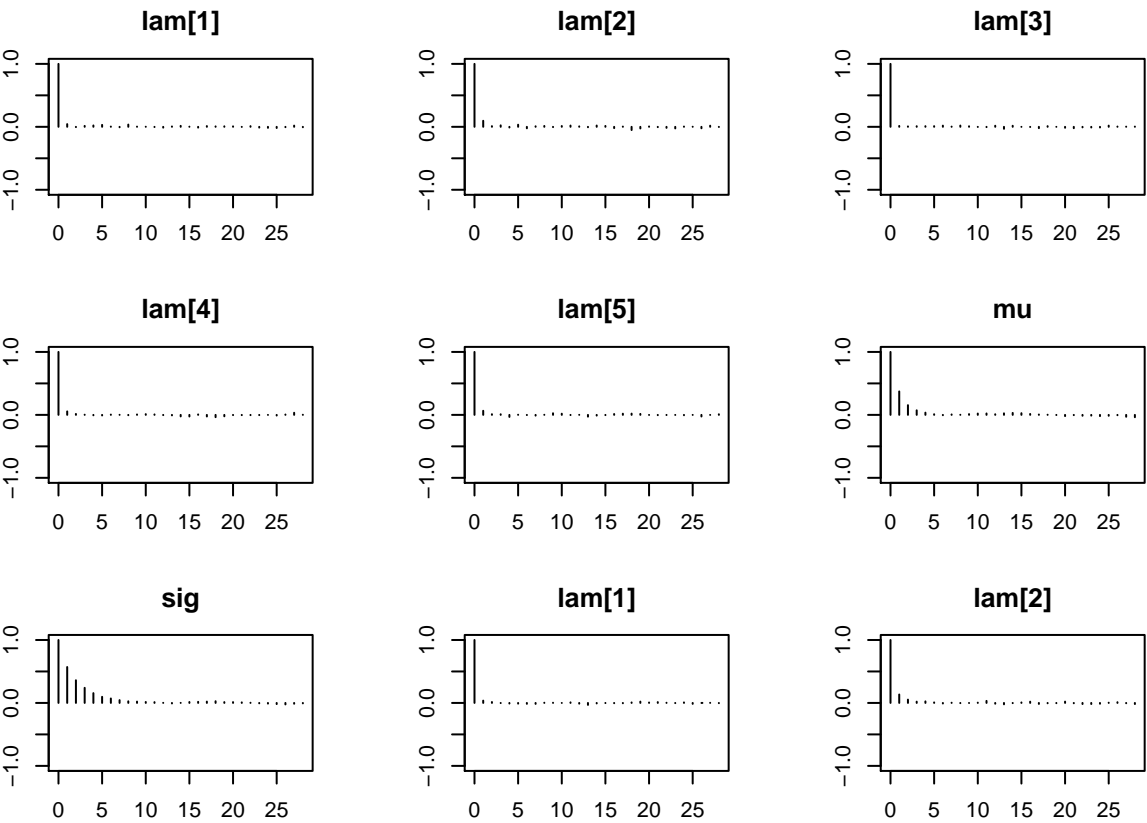
Multivariate psrf

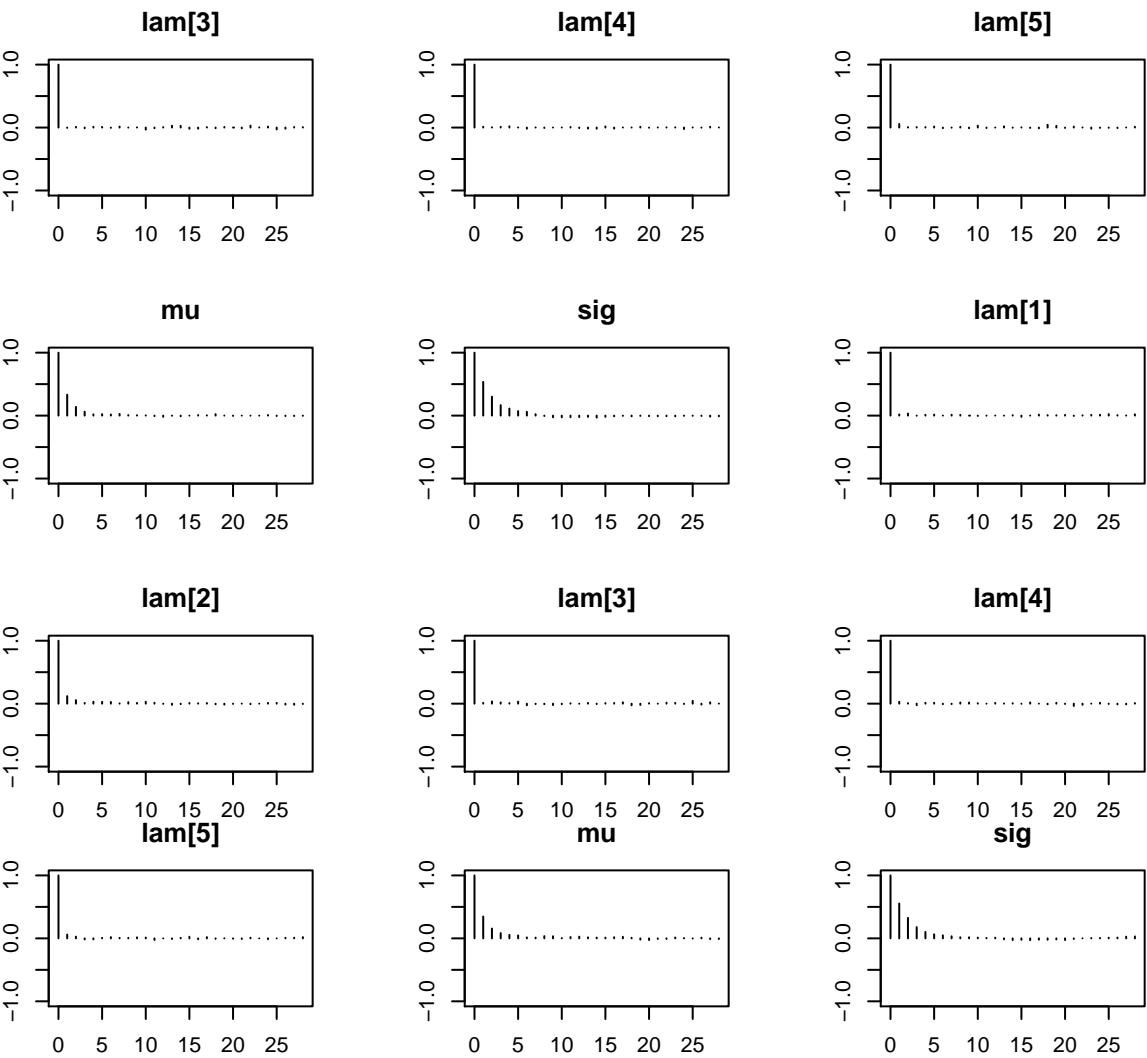
1

	lam[1]	lam[2]	lam[3]	lam[4]	lam[5]
Lag 0	1.000000000	1.000000000	1.000000000	1.000000000	1.000000000
Lag 1	0.032741476	0.11683863	0.006984124	0.032711365	0.060820245
Lag 5	0.012076018	0.02386389	0.017675668	0.002427007	0.008957962



Lag 10	-0.001242761	0.01613687	-0.014641115	0.007045334	0.020754801
Lag 50	0.001173401	0.01369500	-0.001998733	-0.003073964	-0.014490025
	mu		sig		
Lag 0	1.00000000	1.0000000000			
Lag 1	0.35081056	0.5534079582			
Lag 5	0.02705748	0.0776985637			
Lag 10	0.01006484	0.0002367185			
Lag 50	-0.00855037	-0.0070323377			





lam[1]	lam[2]	lam[3]	lam[4]	lam[5]	mu	sig
12998.227	10738.208	14274.512	14186.199	13230.991	6777.891	3999.459

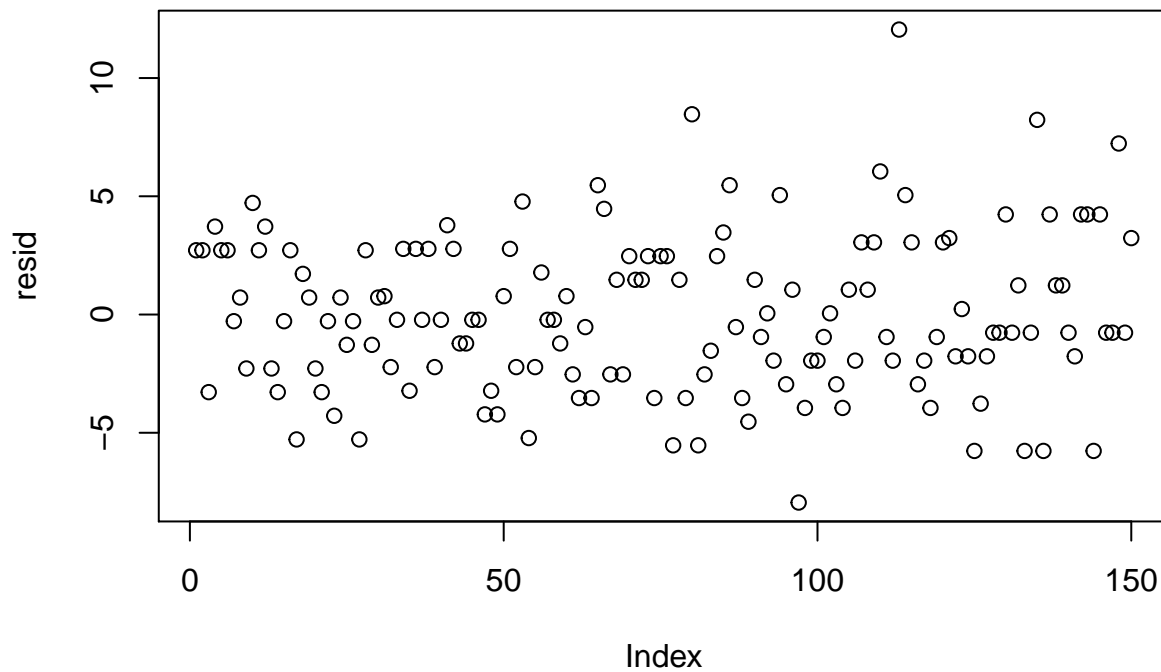
8.1.3 Model checking

After assessing convergence, we can check the fit via residuals. With a hierarhical model, there are now two levels of residuals: the observation level and the location mean level. To simplify, we'll look at the residuals associated with the posterior means of the parameters.

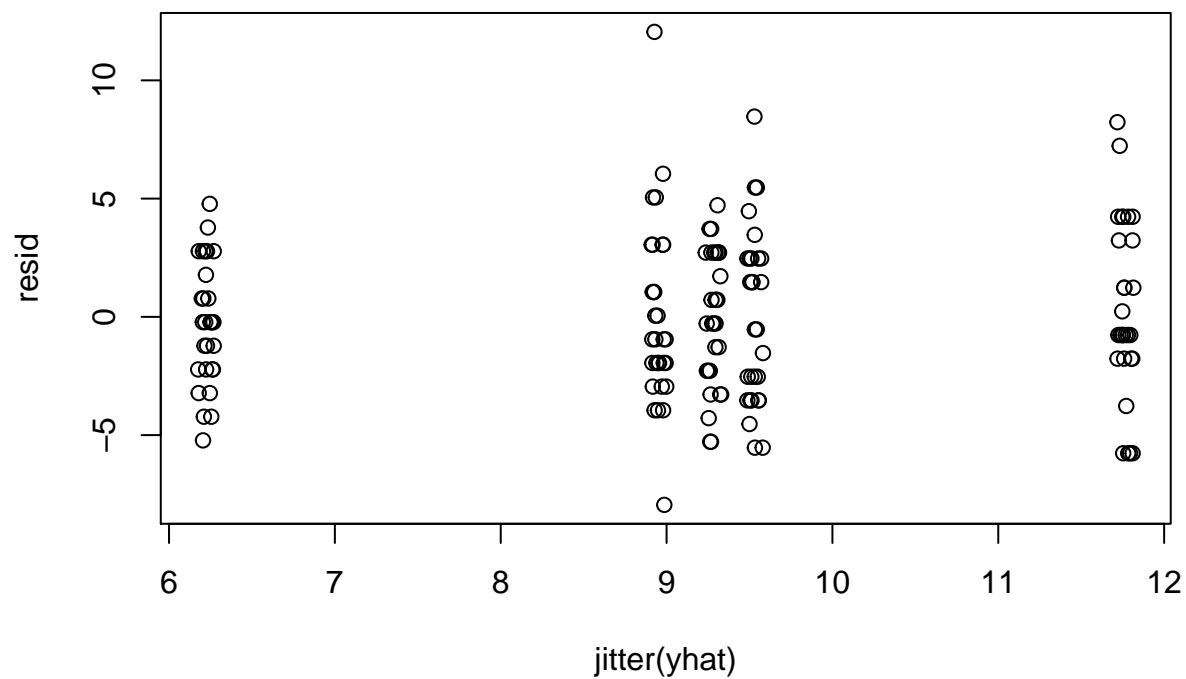
First, we have observation residuals, based on the estimates of location means.

observation level residuals


```
(pm_params = colMeans(mod_csim))  
  
    lam[1]    lam[2]    lam[3]    lam[4]    lam[5]      mu      sig  
9.282030  6.223358  9.531781  8.950109 11.768052  9.111780  2.079838  
  
yhat = rep(pm_params[1:5], each=30)  
resid = dat$chips - yhat  
plot(resid)
```



```
plot(jitter(yhat), resid)
```



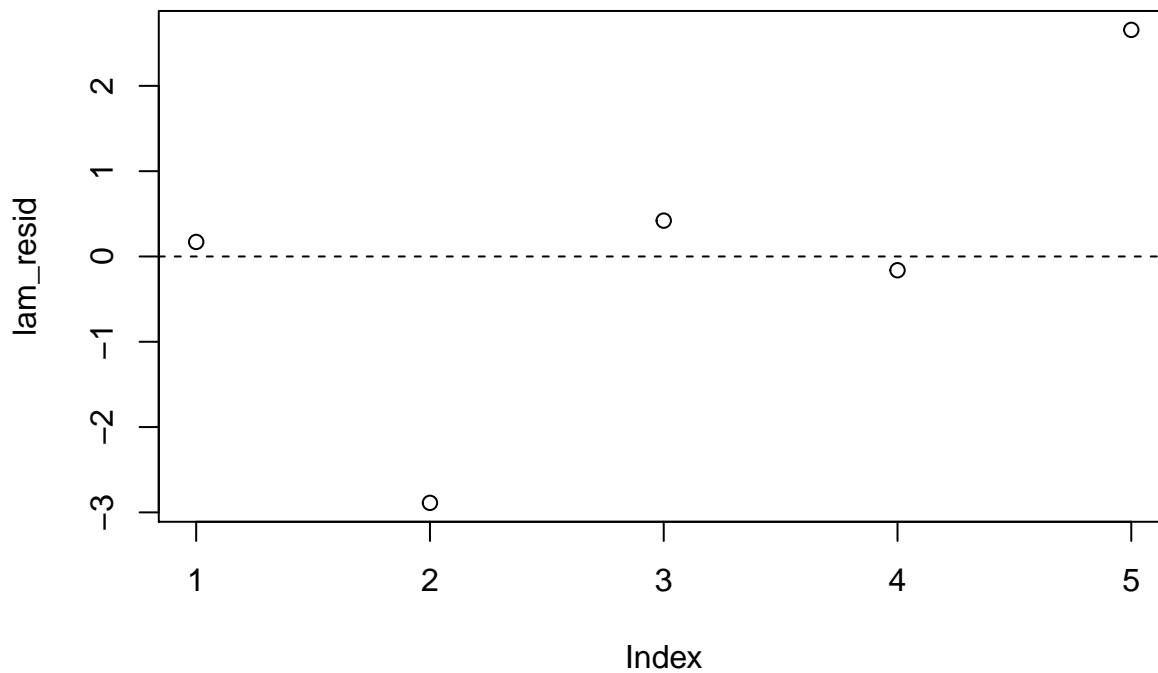
```
var(resid[yhat<7])

[1] 6.447126

var(resid[yhat>11])

[1] 13.72414

## location level residuals
lam_resid = pm_params[1:5] - pm_params["mu"]
plot(lam_resid)
abline(h=0, lty=2)
```



We don't see any obvious violations of our model assumptions.

8.1.4 Results

```
summary(mod_sim)
```

```
Iterations = 2001:7000
```

```
Thinning interval = 1
```

```
Number of chains = 3
```

```
Sample size per chain = 5000
```

1. Empirical mean and standard deviation for each variable,
plus standard error of the mean:

	Mean	SD	Naive SE	Time-series SE
lam[1]	9.282	0.5323	0.004346	0.004687
lam[2]	6.223	0.4694	0.003833	0.004543



lam[3]	9.532	0.5444	0.004445	0.004567
lam[4]	8.950	0.5295	0.004323	0.004450
lam[5]	11.768	0.6173	0.005040	0.005367
mu	9.112	0.9644	0.007874	0.011719
sig	2.080	0.6883	0.005620	0.010955

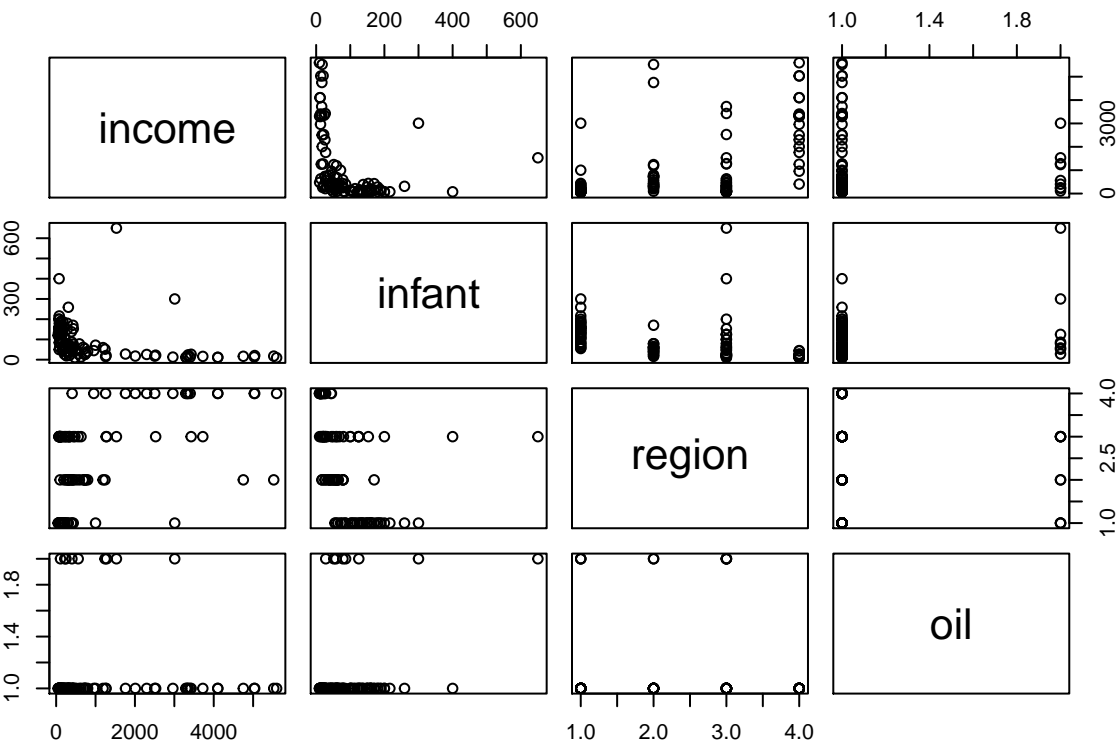
2. Quantiles for each variable:

	2.5%	25%	50%	75%	97.5%
lam[1]	8.273	8.919	9.268	9.628	10.359
lam[2]	5.335	5.903	6.210	6.532	7.162
lam[3]	8.496	9.155	9.522	9.891	10.636
lam[4]	7.949	8.590	8.941	9.301	10.011
lam[5]	10.594	11.350	11.755	12.174	13.027
mu	7.281	8.504	9.077	9.684	11.129
sig	1.105	1.592	1.956	2.431	3.760

8.1.5 Random intercept linear model

We can extend the linear model for the Leinhardt data on infant mortality by incorporating the region variable. We'll do this with a hierarchical model, where each region has its own intercept.

```
'data.frame':  105 obs. of  4 variables:
 $ income: int  3426 3350 3346 4751 5029 3312 3403 5040 2009 2298 ...
 $ infant: num  26.7 23.7 17 16.8 13.5 10.1 12.9 20.4 17.8 25.7 ...
 $ region: Factor w/ 4 levels "Africa","Americas",...: 3 4 4 2 4 4 4 4 4 4 .
 $ oil    : Factor w/ 2 levels "no","yes": 1 1 1 1 1 1 1 1 1 1 ...
```



	income	infant	region	oil
Australia	3426	26.7	Asia	no
Austria	3350	23.7	Europe	no
Belgium	3346	17.0	Europe	no
Canada	4751	16.8	Americas	no
Denmark	5029	13.5	Europe	no
Finland	3312	10.1	Europe	no

Previously, we worked with infant mortality and income on the logarithmic scale. Recall also that we had to remove some missing data.

```
'data.frame':  101 obs. of  6 variables:
 $ income   : int  3426 3350 3346 4751 5029 3312 3403 5040 2009 2298 ...
 $ infant   : num  26.7 23.7 17 16.8 13.5 10.1 12.9 20.4 17.8 25.7 ...
 $ region   : Factor w/ 4 levels "Africa","Americas",...: 3 4 4 2 4 4 4 4 4 4 ...
 $ oil      : Factor w/ 2 levels "no","yes": 1 1 1 1 1 1 1 1 1 1 ...
 $ logincome: num  8.14 8.12 8.12 8.47 8.52 ...
 $ loginfant: num  3.28 3.17 2.83 2.82 2.6 ...
```



```
- attr(*, "na.action")=Class 'omit'  Named int [1:4] 24 83 86 91
.. ..- attr(*, "names")= chr [1:4] "Iran" "Haiti" "Laos" "Nepal"
```

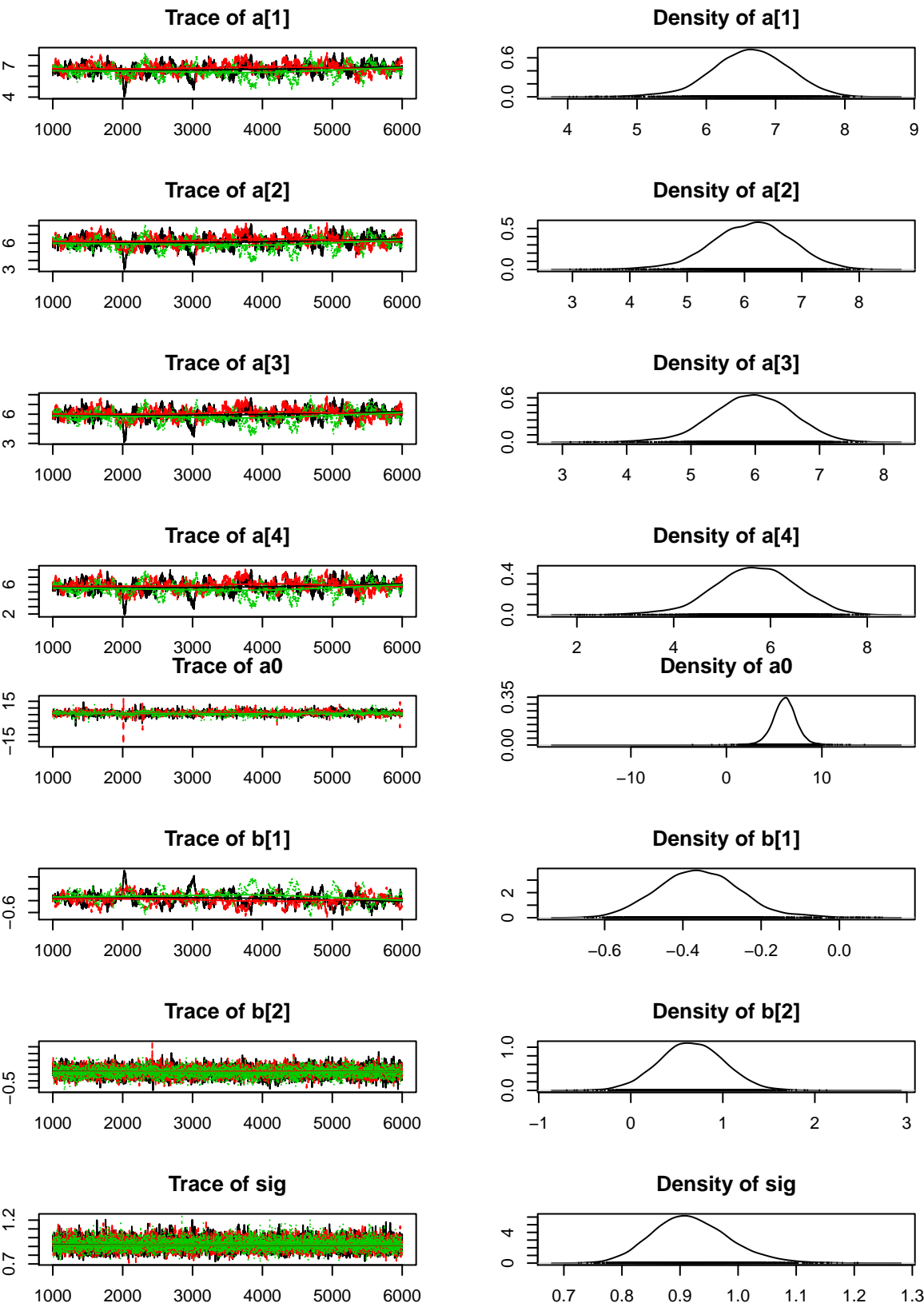
Now we can fit the proposed model:

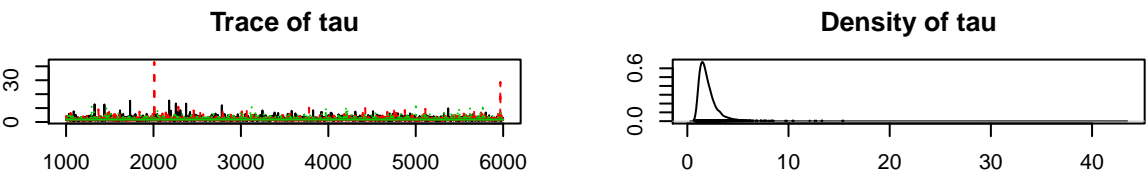
```
[1] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 1 1 1 1 1 1 1 0 0 0 0 0 0 0
[36] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
[71] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
```

```
1 2 3 4
0 31 20 24 18
1 3 2 3 0
```

```
Compiling model graph
  Resolving undeclared variables
  Allocating nodes
Graph information:
  Observed stochastic nodes: 101
  Unobserved stochastic nodes: 9
  Total graph size: 639
```

```
Initializing model
```





Potential scale reduction factors:

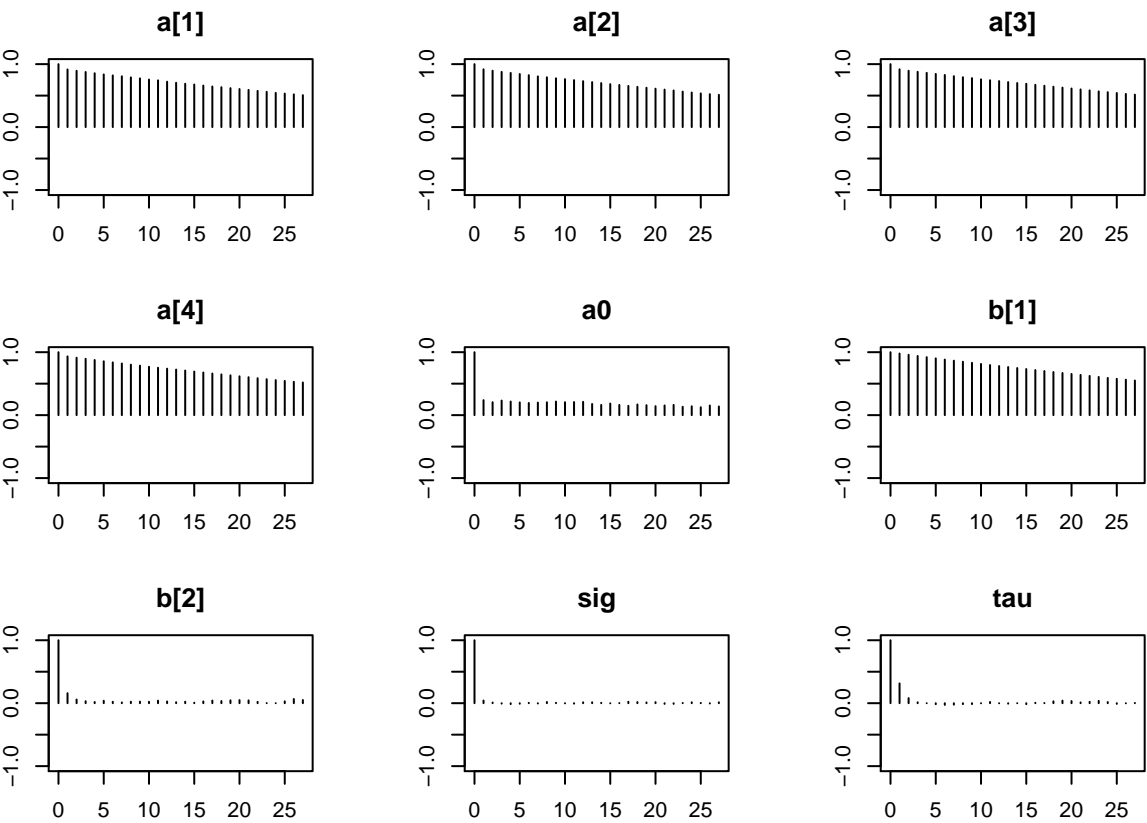
	Point est.	Upper C.I.
a[1]	1.10	1.29
a[2]	1.11	1.31
a[3]	1.10	1.30
a[4]	1.11	1.31
a0	1.02	1.07
b[1]	1.11	1.33
b[2]	1.00	1.01
sig	1.00	1.00
tau	1.01	1.01

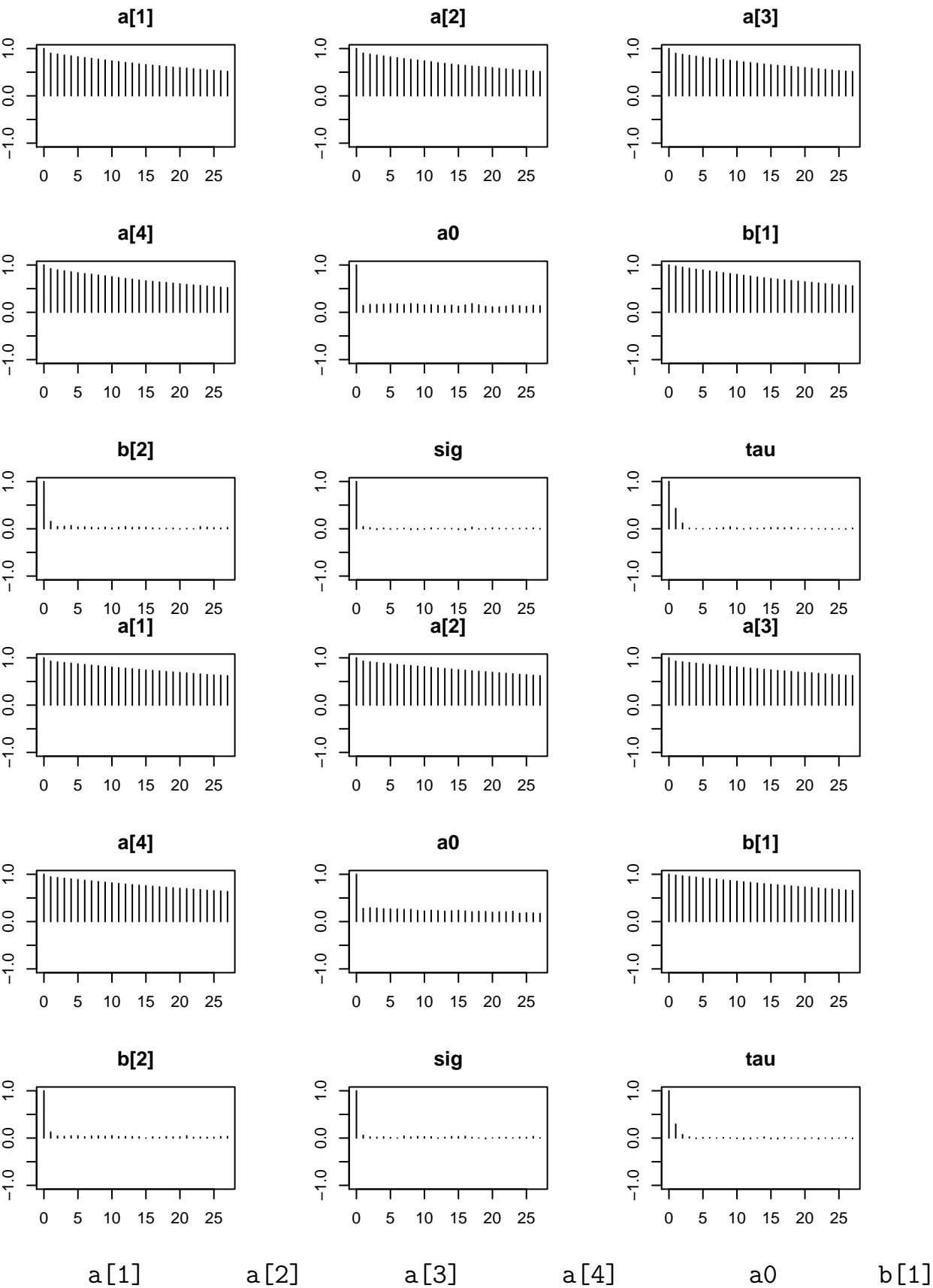
Multivariate psrf

1.07

	a[1]	a[2]	a[3]	a[4]	a0	b[1]
Lag 0	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000	1.0000000
Lag 1	0.9152658	0.9174720	0.9143130	0.9331822	0.21993958	0.9793830
Lag 5	0.8438227	0.8459482	0.8437344	0.8602048	0.21514735	0.9054881
Lag 10	0.7653068	0.7696096	0.7648517	0.7796503	0.19399194	0.8220008
Lag 50	0.3605857	0.3592238	0.3590416	0.3645316	0.09356611	0.3846434
	b[2]	sig	tau			
Lag 0	1.000000000	1.000000000	1.000000000			
Lag 1	0.147024106	0.048178001	0.3463555333			
Lag 5	0.040588329	-0.001527570	-0.0029154546			

Lag 10 0.028718120 0.005505842 0.0028262763
Lag 50 0.008303821 -0.002587573 -0.0007569181







162.0426	152.8432	164.8619	159.3883	709.9542	146.6284
b[2]	sig	tau			
6345.4749	11725.2325	7787.5196			



8.2 Meta analysis



Bibliography

Efron, B., 2013. Bayes theorem in the 21st century. *Science* 340, 1177–1178.