

# A simple example of Bayesian statistics and MCMC estimation

Michael Ash\*

20 June 2019

The goal is to estimate the parameters (mean  $\mu$  and variance  $\sigma^2$ ) of normally distributed data,

$$y_i \sim N(\mu, \sigma^2), \quad i = 1, \dots, n$$

or more accurately to estimate the posterior distribution of the parameters.

The mean is  $\mu$ , but instead of  $\sigma^2$  we will examine  $h$ , the inverse of the variance, i.e.,  $h = 1/\sigma^2$ . Hence the actual target is

$$y_i \sim N(\mu, h^{-1}), \quad i = 1, \dots, n.$$

The reason for this construction is that it's more convenient to use a Gamma distribution for the distribution of  $h$  rather than an inverted- (or inverse-) Gamma distribution for the distribution of  $\sigma^2$ .

One of the key pieces of Bayesian statistics is the prior distribution, subjective belief about the distribution of the parameters before we examine the data. There are two parameters,  $\mu$  and  $h$ , for which we must have prior beliefs.

The prior for  $\mu$  is

$$\mu \sim N(\mu_0, h_0^{-1}).$$

Thus, our belief about  $\mu$  is that it is normally distributed with mean  $\mu_0$  and variance  $h_0^{-1}$ . The density function for this prior, but for some proportionality terms, is  $\pi(\mu) \propto \exp[-\frac{h_0}{2}(\mu - \mu_0)^2]$ . Note that this is a rather diffuse prior distribution because the guess about the spread of our beliefs about  $\mu$ , the mean of  $y$ , is the same as the  $\sigma^2$  variance of the underlying distribution of  $y_i$  itself (not, for example, the CLT-implied  $\sigma^2/n$  which is a much tighter distribution—note that we don't yet have data, including  $n$ ).<sup>1</sup>

The prior for  $h$  is

$$h \sim G(\alpha_0/2, \delta_0/2),$$

where  $G(\cdot, \cdot)$  is the Gamma distribution. The Gamma distribution depends on two parameters (“hyperparameters” with respect to the main model), in this

---

\*The main source for this document, including notation, is Edward Greenberg, *Introduction to Bayesian Econometrics*. Cambridge University Press, 2008.

<sup>1</sup>In the context of the distribution of  $\mu$ ,  $h$  is referred to as “precision,” the inverse of “variance”.

case,  $\alpha$  and  $\delta$  with initial (“prior”) values of  $\alpha_0$  and  $\delta_0$ . Although the Gamma distribution is somewhat obscure, it’s probably easiest to think about its parameters by noting a special case of the Gamma distribution:  $G(\nu/2, 1/2)$  is the same as a chi-squared distribution  $\chi_\nu^2$ . First, this connects the reader to the inverted chi-squared distribution of estimates of the variance in frequentist statistics. As with the chi-squared distribution and as to be expected with variance, the density is only non-zero for positive arguments, i.e., the variance must be positive. Second the special case gives some sense of the magnitude of the metaparameters,  $\alpha$  around twice the sample size and  $\delta$  around twice the sum of squared errors are plausible values (note that both terms get halved as arguments in  $G(\cdot, \cdot)$ ).<sup>2</sup> The density for this prior, again but for some proportionality terms, is  $\pi(h) \propto h^{\alpha_0/2-1} \exp\left[-\frac{\delta_0 h}{2}\right]$ .

It’s worth pointing out at this point that the prior  $\mu_0$  for  $\mu$  depends on the prior  $h_0$  for  $h$  but not vice versa. The prior density for both that we work with will really be

$$\pi(\mu, h) = \pi(\mu|h)\pi(h).$$

The object of interest in Bayesian statistics is the posterior distribution, which unites the prior with the data to produce updated beliefs about the distribution of the parameters. We’ve discussed the priors above. The data are represented by the likelihood function, which expresses how likely are these data if the parameters were given for the random process that generated the data,  $f(y|\mu, h)$ . (How likely are the data  $y$  for given values of the parameters of a normal distribution,  $\mu$  and  $h$ ?)

Note that the posterior, the object of interest, is how likely are (each of the possible values of) the parameters with the data given, or  $\pi(\mu, h|y)$ . Bayes formula is what allows us to compute the posterior, the object of interest, given the data and the priors:

$$\pi(\mu, h|y) \propto f(y|\mu, h) \times \pi(\mu, h).$$

It’s helpful to read this out loud: After seeing the data, how likely is a particular value of  $\mu$ ? It depends on how likely that value of  $\mu$  was initially considered to be times how likely the observed values of the data were to occur if that value of  $\mu$  were in fact the case.

Bayes formula also demonstrates why we need priors; it would otherwise be impossible to convert the likelihood function, where the probability of the *data given parameters*, is the output, into the posterior distribution, where the probability of *parameters given data*, is the output.

By the way, the formula uses “ $\propto$ ” (“is proportional to”) rather than “=” because a normalization term is dropped from the denominator of Bayes’s formula.<sup>3</sup>

---

<sup>2</sup>The parameters are  $G(\text{shape}, \text{rate})$  with “rate” sometimes described as “inverse scale.” In many contexts, shape and rate are denoted  $\alpha$  and  $\beta$ , which we’ll avoid outside this footnote because  $\alpha$  is otherwise used. For  $X \sim G(\alpha, \beta)$ ,  $E(X) = \alpha/\beta$  and  $V(x) = \alpha/\beta^2$ . In terms of the  $\alpha$  and  $\delta$  in the main body of the paper,  $E(h) = \alpha/\delta$  and  $V(h) = \alpha/(\delta^2/4)$ . I made up guesses for the priors on  $h$  by checking that these terms gave plausible values.

<sup>3</sup>“But for a proportionality term” appears all over the Bayesian literature. The actual

Since we know that the data were generated from a normal distribution, we know the functional form of the likelihood function, which, again, expresses how likely is this particular outcome  $(y_1, \dots, y_n)$  with given parameters  $\mu$  and  $h$ . In this case,

$$f(y|\mu, h) = h^{n/2} \exp \left[ -\frac{h}{2} \sum (y_i - \mu)^2 \right].$$

We have already specified the priors for the two parameters. We can multiply the likelihood by the priors to derive the posterior:

$$\pi(\mu, h|y) \propto \overbrace{h^{n/2} \exp \left[ -\frac{h}{2} \sum (y_i - \mu)^2 \right]}^{\text{Likelihood function}} \times \overbrace{\exp \left[ -\frac{h_0}{2} (\mu - \mu_0)^2 \right]}^{\text{Prior density of } \mu} \times \overbrace{h^{\alpha_0/2-1} \exp \left[ -\frac{\delta_0 h}{2} \right]}^{\text{Prior density of } h}.$$

The first term on the right-hand side is the likelihood function and it is the only place where the data, the  $y_i$ , appear. The second term is the prior normal distribution of  $\mu$  and the third term is the prior Gamma distribution of  $h$ . The analysis takes advantage of the dependence of  $\mu$  on  $h$  and of the independence of  $h$ ; so the prior terms can be simply multiplied together.

At this point, we additively combine the exponential terms, rearrange, add like terms, notice some instances of the mean of the data ( $\bar{y} = \sum y_i/n$  is the MLE of  $\mu$ ), etc.<sup>4</sup> The rearranged equation can be elegantly split into conditional posterior distributions of  $h$  and  $\mu$ .

The conditional posterior density of  $h$  is

$$\pi(h|\mu, y) \propto h^{(\alpha_0+n)/2-1} \exp \left[ -h \frac{\delta_0 + \sum (y_i - \mu)^2}{2} \right],$$

which is the density function of  $G[(\alpha_0+n)/2, (\delta_0 + \sum (y_i - \mu)^2)/2]$ . Thus, we have updated parameters for the Gamma distribution which describe the posterior Gamma distribution of  $h$ .

The old, prior first parameter was  $\alpha_0/2$ . The new, data-incorporating first parameter is  $(\alpha_0 + n)/2$ . The old, prior second parameter was  $\delta_0/2$ . The new data-incorporating second parameter is  $(\delta_0 + \sum (y_i - \mu)^2)/2$ .

---

expression is  $\pi(\mu, h|y) = \frac{f(y|\mu, h) \times \pi(\mu, h)}{f(y)}$ . The  $f(y)$  in the denominator is necessitated by Bayes's formula  $P(A \cap B) = P(A|B)P(B) = P(B|A)P(A)$ , i.e.,  $f(y)$  is the  $P(B)$  term, which is needed to normalize the expression for  $P(A|B) = P(B|A)P(A)/P(B)$  to be between 0 and 1. It is ok to replace “=” with “ $\propto$ ” thereby “ignoring”  $f(y)$  because the marginal distribution of the data does not vary and, hence,  $1/f(y)$  simply scales the density distribution up or down without altering the location of peaks or valleys or areas of portions relative to the area of the whole. This replacement of “=” with “ $\propto$ ” is particularly useful with continuously distributed parameters where relative probabilities are of interest. If the question involves parameters for a discrete distribution, e.g., the probability that a coin is fair, it may be worthwhile to carry through the normalization term and get actual posterior probabilities with “=”.

<sup>4</sup>This rearrangement is a bit tricky. It involves recognizing the parameters of the posterior density function (in  $\mu$  and  $h$ ) as constructed from summary statistics of the data and parameters of the prior. At one point the rearrangement involves “completing the square.”

The conditional posterior density of  $\mu$  is

$$\pi(\mu|h, y) \propto \exp \left[ -\frac{h_0 + hn}{2} \left( \mu - \frac{h_0\mu_0 + hn\bar{y}}{h_0 + hn} \right)^2 \right],$$

which can be recognized as the density formula for a variable distributed  $N[(h_0\mu_0 + hn\bar{y})/(h_0 + hn), (h_0 + hn)^{-1}]$ . That is, the posterior distribution of  $\mu$  is distributed normally with  $E(\mu) = (h_0\mu_0 + hn\bar{y})/(h_0 + hn)$  and  $\sigma_\mu^2 = (h_0 + hn)^{-1}$ .

A reason to choose a normal distribution for the prior of  $\mu$  and a Gamma distribution for the prior of  $h$  is that, for the parameters of a normal distribution (which define the likelihood function), these priors generate a normal distribution for the posterior of  $\mu$  and a Gamma distribution for the posterior of  $h$ . When the posterior distribution uses the same functional form (with new parameter values, of course) as the prior, they are said to be “conjugate” and the problem is more tractable.

In this example it is reasonably easy to trace how the data, initially embedded only in the likelihood function, and the prior combine to form the posterior. We can examine the expected value of the mean  $\mu$  in more detail:

$$E(\mu) = (h_0\mu_0 + hn\bar{y})/(h_0 + hn)$$

and rewrite this term to highlight how the posterior  $E(\mu)$  is a weighted average<sup>5</sup> of the prior estimate of the mean,  $\mu_0$ , and the empirical (maximum likelihood) estimate of the mean,  $\bar{y}$ :

$$E(\mu) = \frac{h_0}{h_0 + hn} \mu_0 + \frac{hn}{h_0 + hn} \bar{y}.$$

Note in particular how the relative weight on the second term depends on the sample size  $n$ . As the sample size gets larger and larger, the second weight approaches 1 (and the first weight approaches 0); hence the empirical estimate of the mean increasingly dominates the prior estimate as the sample gets large.<sup>6</sup>

Note also that the estimate of the posterior spread of beliefs about  $\mu$  is  $\sigma_\mu^2 = (h_0 + hn)^{-1} = \frac{1}{1/\sigma_{y0}^2 + n/\sigma_y^2}$ , which looks something like the frequentist CLT estimate  $\sigma_\mu^2 = \sigma_y^2/n$  (modified slightly by the inclusion of the prior estimate  $h_0 = 1/\sigma_{y0}^2$ ).

In some sense, we are now done because we have written down a concise analytic expression for the posterior. But suppose that the posterior was very messy, and we were unable to write down an analytic expression. For example,

---


<sup>5</sup>Recall that the weights in a weighted average sum to 1 (in this case  $\frac{h_0}{h_0 + hn} + \frac{hn}{h_0 + hn} = 1$ ) and express how much importance each of the terms being averaged carries in the average.

<sup>6</sup>I am less familiar with the Gamma distribution but the terms can be similarly interpreted to see the growing importance of the empirical estimate of the variance relative to the prior distribution.

the two parameters might not be independently distributed or the combination of the prior with the likelihood might be messy. In this case, we can use simulation to generate a posterior distribution for the parameters.<sup>7</sup>

Below is a Gibbs algorithm that randomly draws a new value of each parameter in successive steps. Note the iteration between the previous value of one parameter and the other. Some key elements of the simulation approach are the use of one previous value and the iterated conditionality, i.e., the next draw of one parameter depends on and only on one previous draw of the other parameter.

1. Choose a starting value for  $\mu = \mu^{(0)}$ .<sup>8</sup>
2. Sample  $h^{(1)}$  from  $G(\alpha_1/2, \delta_1/2)$ , where  $\alpha_1 = \alpha_0 + n$  and  $\delta_1 = \delta_0 + \sum (y_i - \mu^{(0)})^2$ .
3. At the  $g$ th iteration, draw
  - $\mu^{(g)}$  from  $N[(h_0\mu_0 + h^{(g-1)}n\bar{y})/(h_0 + h^{(g-1)}n), (h_0 + h^{(g-1)}n)^{-1}]$
  - $h^{(g)}$  from  $G[\alpha_1/2, (\delta_0 + \sum (y_i - \mu^{(g)})^2)/2]$

After discarding a number of burn-in draws, store the  $\mu^{(g)}$  and  $h^{(g)}$ . These constitute the simulated distributions of  $\mu$  and  $h$  and can be analyzed for means, credibility intervals (analogous to confidence intervals in frequentist statistics), etc.  code and plots of the prior and posterior distributions follow.

```

1  ## Run with R CMD BATCH --vanilla gibbs-normal-parameters.R &
2
3  ## Birth weight data from UC Berkeley
4  babiesI <- read.table("http://www.stat.berkeley.edu/users/statlabs/data/babiesI.data", header=TRUE)
5  ## babiesI <- read.table("babiesI.data", header=TRUE)
6
7  ## 2000 burn-in and 5000 post-burn-in iterations
8  burnin <- 2000
9  h <- rep(NA, 7000)
10 mu <- rep(NA, 7000)
11 delta <- rep(NA, 7000)
12
13 ## The truth
14 with(babiesI, mean(bwt))
15 with(babiesI, 1/var(bwt))
16 with(babiesI, sd(bwt))
17
18 ## Draw a sample from the population
19 bwt <- with(babiesI, sample(bwt, size=100, replace = FALSE))
20
21 (n <- length(bwt))
22 (ybar <- mean(bwt))
23 sd(bwt)
24
25 ## Priors
26 mu0 <- 110
27 alpha0 <- 30
28 delta0 <- 6000
29 (h0 <- (alpha0/2) / (delta0/2))

```

<sup>7</sup>Note that we got the parameters by analytic means, not estimation by simulation, which is a different method. It's the distribution of the posterior which we will ostentatiously get by simulation using the analytically derived parameters.

<sup>8</sup>If desired the sampling can begin with  $h^{(0)}$  and the algorithm modified accordingly.

```

30 ( 1 / h0 )
31 sqrt ( 1 / h0 )
32
33 alpha1 <- alpha0 + n
34 delta[i] <- delta0 + sum((bwt - mu0)^2)
35 h[i] <- rgamma(n=1,shape=alpha1/2,rate=delta[i]/2)
36 mu[i] <- rnorm(n=1,mean = (h0*mu0 + h0*n*ybar) / (h0 + h0*n), sd=sqrt ( 1 / (h0 + h0*n) ) )
37
38 ## Here are the draws from the simulated distribution
39 for(i in 2:length(mu)) {
40   mu[i] <- rnorm(n=1,
41     mean = (h0*mu0 + h[i-1]*n*ybar) / (h0 + h[i-1]*n),
42     sd=sqrt ( 1 / (h0 + h[i-1] * n) ) )
43   delta[i] = (delta0 + sum((bwt - mu[i])^2))
44   h[i] <- rgamma(n=1,shape=alpha1/2,rate= delta[i] / 2)
45 }
46
47 ## Report the mean and the 95 percent credibility interval of the
48 ## posterior of mu
49 mean(mu[(burnin+1):length(mu)])
50 quantile(mu[(burnin+1):length(mu)], probs=c(0.025,0.975))
51
52 ## Report the mean and 95 percent c.i. of the posterior of the
53 ## precision, the variance, and the standard deviation
54 mean(h[(burnin+1):length(h)])
55 quantile(h[(burnin+1):length(h)], probs=c(0.025,0.975))
56
57 (1/mean(h[(burnin+1):length(h)]))
58 quantile((1/h[(burnin+1):length(h)]), probs=c(0.025,0.975))
59
60 sqrt(1/mean(h[(burnin+1):length(h)]))
61 quantile(sqrt(1/h[(burnin+1):length(h)]), probs=c(0.025,0.975))
62
63 ## Prior of mean (mu)
64 pi.mu0 <- dnorm(x.mu <- seq(from=70,to=150,by=(150-70) / (length(mu)-burnin) )
65   , mean=mu0, sd=sqrt(1/h0))
66 ## Posterior of mean (mu) with prior overlay
67 hist(mu[(burnin+1):length(mu)],xlim=c(70,150),main="",xlab="")
68 lines(x.mu,30000*pi.mu0)
69 title(expression(paste("Prior and Posterior Distribution of ", mu)),xlab=expression(mu))
70
71 ## Prior of precision (h)
72 pi.h0 <- dgamma(x.h <- seq(from=(1/30)^2,to=(1/10)^2,by=((1/10)^2-(1/30)^2)/ (length(mu)-burnin) ),
73   shape=alpha0 / 2, rate=delta0 / 2)
74 ## Posterior of precision (h) with prior overlay
75 hist(h[(burnin+1):length(h)],xlim=c((1/30)^2,(1/10)^2),main="",xlab="")
76 lines(x.h,2.7*pi.h0)
77 title(expression(paste("Prior and Posterior Distribution of ", h)),xlab=expression(paste("h")))
78
79 ## Prior of standard deviation (sqrt(1/h))
80 ## Posterior of standard deviation (sqrt(1/h)) with prior overlay
81 hist(sqrt(1/h[(burnin+1):length(h)]),xlim=c(10,30),main="",xlab="")
82 lines(sqrt(1/x.h),2.5*pi.h0)
83 title(expression(paste("Prior and Posterior Distribution of ", sigma)),xlab=expression(sigma))

```

