

# Bayes Computational Methods

Evan Miyakawa

## Introduction

### The Bayes Paradigm

In Bayesian statistical inference, we are interested in finding the posterior distribution of a certain parameter, given a prior distribution and a likelihood for that parameter from a sample of data. Bayes Rule states that, given a parameter of interest  $\theta$  and some known data  $x$ , we can find the conditional probability of  $\theta$  given  $y$  using

$$p(\theta|x) = \frac{p(x|\theta)p(\theta)}{p(x)},$$

where in the discrete case  $p(x) = \sum_{\theta} p(\theta)p(x|\theta)$ , and in the continuous case,  $p(x) = \int p(\theta)p(y|\theta)d\theta$ . More formally, we can derive the posterior distribution of a parameter  $\theta$ ,  $p(\theta|x)$ , by

$$p(\theta|x_1, x_2, \dots, x_n) = \frac{\ell(\theta|x_1, x_2, \dots, x_n)\pi(\theta)}{\int \ell(\theta|x_1, x_2, \dots, x_n)\pi(\theta)d\theta},$$

where  $\ell(x|\theta)$  is proportional to the likelihood of  $\theta$ , conditioned on the data, and  $\pi(\theta)$  is the prior distribution of  $\theta$ . For example, suppose that  $X_i \sim \text{Bern}(\theta)$ , and we are interested in inference on  $\theta$ . A reasonable prior to assign to  $\theta$  would be a  $\text{Beta}(\alpha, \beta)$  distribution, given that  $\theta$  is constrained to be between 0 and 1, which is a property of the Beta distribution. Therefore, given that  $\theta \sim \text{Beta}(\alpha, \beta)$  for some known values of  $\alpha$  and  $\beta$ , then

$$\begin{aligned} p(\theta|x) &= \frac{\ell(\theta|x_1, x_2, \dots, x_n)\pi(\theta)}{\int \ell(\theta|x_1, x_2, \dots, x_n)\pi(\theta)d\theta} \\ &= \frac{\theta^{\sum_{i=1}^n x_i} (1-\theta)^{n-\sum_{i=1}^n x_i} \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1}}{\int \theta^{\sum_{i=1}^n x_i} (1-\theta)^{n-\sum_{i=1}^n x_i} \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1} d\theta} \\ &\propto \theta^{\sum(x_i)+\alpha-1} (1-\theta)^{n-\sum(x_i)+\beta-1}, \end{aligned}$$

which is the kernel of a  $\text{Beta}(\sum(x_i) + \alpha, n - \sum(x_i) + \beta)$  distribution. Therefore, we can analytically calculate the posterior distribution of  $\theta$  as

$$p(\theta|x) \sim \text{Beta}\left(\sum(x_i) + \alpha, n - \sum(x_i) + \beta\right).$$

### Bayesian Computation

In more complicated scenarios, it is not easy to calculate the posterior distribution of a parameter by hand, and in many cases it is impossible. This is where computing can help us. If we could somehow computationally grab samples from the posterior distribution we are interested in, then we could obtain most of the statistics that we are interested in.

## History of Bayesian computation up to MCMC

### Markov Chain Monte Carlo

Markov Chain Monte Carlo (referred to as MCMC), is a common approach that allows for drawing samples of a parameter from a distribution that is approximately equivalent to the target posterior distribution of that parameter. It is an iterative algorithm, where each realization gets closer to the target density we are interested in. An essential property of any Markov chain is that the distribution of each draw only depends on the most recent value drawn before it. As we keep obtaining draws, each draw is probabilistically more likely to have come from the target posterior distribution. Based on the convergence properties of MCMC, we can expect each draw to be coming from the real distribution of the parameter we are interested, once we have let the Markov chain run long enough.

- MCMC is the only suitable method when we don't have a conjugate prior.

### History up through BUGS (include slice sampling)

- BUGS project began in 1989
- BUGS language contains for a somewhat concise list of functions and distributions used for specifying relationships between
- The BUGS language only has a somewhat limited syntax, but is able to describe extremely complex models using this syntax.
- While the programs that use BUGS-type language are continually evolving, the syntax has stayed very stable.
- BUGS language is declarative, which means the order of statements doesn't matter, as opposed to a procedural language.
- A main initial purpose for the creation of bugs was to allow for graphical models – models that can be specified through click and point software that creates nodes, or expressing graphics of models that already exist for easier understanding.
- Six steps in running BUGS models
  1. Checking syntax of model specification
  2. Loading data
  3. Compiling the BUGS model
  4. Starting the simulation at appropriate initial values
  5. Simulate unknown parameters
  6. Gather results and report summary statistics, along with other appropriate output
- WinBUGS was the first widely used program to utilize the BUGS language. However, OpenBUGS is the more recently maintained version, which communicates better with other software and is more portable.
- R2OpenBUGS is the R package that allows for R to interface with OpenBUGS
- 

### Gibbs Sampler

- One of most common algorithms for simulating MCMCs.
- Special case of Metropolis-Hastings alg
- Generates a multi-dimensional markov chain by splitting the vector of random unknown variables and sampling each subvector one by one, conditional on the most recent values of the other variables.

- Algorithm: If we have a vector of unknown parameters  $\theta$ , where  $\theta = (\theta_1, \theta_2, \dots, \theta_k)$ ,
  1. Choose arbitrary starting values  $\theta_1^{(0)}, \theta_2^{(0)}, \dots, \theta_k^{(0)}$
  2. Complete one iteration of the sampler by getting new values for  $\theta$  using the following:
    - Sample a new value for  $\theta_1$ , from the full conditional distribution of  $\theta_1$ , given the most recent values of all other elements of  $\theta$  and the data:

$$\theta_1^{(1)} \sim p(\theta_1 | \theta_2^{(0)}, \theta_3^{(0)}, \dots, \theta_k^{(0)}, y).$$

- Sample a new value  $\theta_2^{(1)}$  for the second component of  $\theta$ , from the full conditional distribution  $p(\theta_2 | \theta_1^{(1)}, \theta_3^{(0)}, \dots, \theta_k^{(0)}, y)$ .
  - Continue this process for the remainder of the  $k$  elements of  $\theta$ .
- 3. Repeat step 2 for the total decided number of iterations that the sampler will run.
- The Gibbs sampling strategy often allows for sampling from high-dimensional spaces to be simplified to drawing from one-dimensional distributions.
- Weakness: because it uses a random-walk behavior, if two variables are highly correlated, traversing the space can take a very long time.

## Metropolis

### HMC

Radford-Neal HMC Chapter: - Hamiltonian Monte Carlo uses Hamiltonian Dynamics - Visualize the dynamics as a frictionless puck that slides over a surface of varying height - The state of the system has three elements: position, momentum, and potential energy - The position corresponds to variables of interest - The potential energy will be minus the log of the probability density for these variables - There is one momentum variable for each position variable - Using “Hamilton’s equations”, the movement of the puck is computed. - For computational implementation, these equations must be approximated by discretizing time, using some small stepsize. - To use Hamiltonian dynamics to sample from a distribution using MCMC, we need to translate the density function for the distribution to a potential energy function. The variables of interest are now position variables - We simulate a markov chain where each iteration resamples the momentum variable associated with each position variable and then does a Metropolis update with the proposal found using Hamiltonian dynamics - Two steps for each iteration: 1. New values for the momentum variables are randomly drawn from their Gaussian distribution, independent of the current values of the position variables. 2. A metropolis update is performed, using hamiltonian dynamics to propose a new state. The proposed state is either accepted or rejected based on the HMC probability function used. If the proposed state is rejected, the next state is the same as the current one. - The HMC algorithm typically will not be trapped in some subset of the state space, which allows for it to always asymptotically converge to the target distribution, and quickly. - Tuning Parameters for HMC: - Stepsize - the leapfrog stepsize. Too large, and it will lead to a low acceptance rate for proposed states. Too small and it will be slow computationally. - Trajectory length -

### STAN

STAN 2015 - Based on a probabilistic programming language that is more flexible and expressive than BUGS/JAGS. - Based on HMC - Works in command line, Python, and R

## Bayesian Samplers

### OpenBUGS

OpenBugs is the oldest of the Bayesian samplers that we use. OpenBugs is designed for Windows machines, but there is a workaround to run it on MacOS and Linux, using WINE, a program designed to run Windows

software on other operating systems. OpenBugs can be used through the actual click-and-point application, or through R, using the `R2OpenBugs` package, which allows for the user to send model scripts and data to OpenBugs without needing to leave the R interface. The R command to run OpenBugs is `bugs()`. Using this method causes OpenBugs to run invisibly in the background and return the results, unless an argument is specified to debug code in the actual application itself, which can be very useful. When using R to run OpenBugs, there is little way to know what went wrong if the code errors somewhere without interacting with the OpenBugs program itself. Using the `debug = T` argument within `bugs()` will cause the application to run in the forefront as if it were actually being used directly by the user. This way, any messages given by the application while running can be viewed directly. It is recommend to use `debug = T` only when debugging, as it does slow down the sampling process slightly. OpenBugs has many tools to visualize the performance the sampler which are easy to use. The biggest drawback of using the OpenBugs application directly is that specifying the data can be tedious if the data is not very small in size or design. Using R to specify the data is much easier. One unique feature in OpenBugs is the Doodle tool, which allows for model specification through the creation of a diagram of nodes representing the parameters. Though using Doodle can be helpful for illustration, it is not used often in practice.

## JAGS

In terms of model specification, JAGS is very similar to OpenBugs, but works more seamlessly in R using the package `rjags`. JAGS was developed independently of OpenBUGS, but uses the BUGS language for most of its model specification syntax. The `runjags` is an optional package that allows for the JAGS model script to be inputted as a string, and run using the `run.jags()` function. Since JAGS runs inside R itself, it is very easy to work with. Helpful errors are given when there is an issue with the model code or sampling process. JAGS is often a very fast sampler, especially for less complicated models.

## Nimble

Nimble is an extension of Bugs that allows for bayesian computation within R, but uses C++ under the hood to compile models and algorithms for speed. Using the `nimble` package, Nimble can perform MCMC bayesian computations, but also allows for other types of algorithms to be loaded or custom-built. It can also compile other code without BUGS models into C++ for other uses. For the purpose of bayesian sampling using MCMC, Nimble is used very similarly to BUGS in practicality. The flexibility of Nimble is not needed in most cases when doing MCMC. The model specification language is almost identical to OpenBugs. A model can be both compiled and run using the function `nimbleMCMC()`. Because the model is compiled using C++ each time, this can drastically increase the total computation time as compared to other sampling methods. Occasionally, Nimble can also freeze up or fail, which can take a bit of time to fix.

## STAN

STAN is a language in-and-of itself used for statistical modeling and data analysis. STAN can be used in several programming languages, including R and Python. While many bayesian samplers use Gibbs sampling, STAN uses Hamiltonian Monte Carlo under the hood, a more computationally efficient sampler that is based on an advanced physics simulation. This can greatly reduce the computational time for many types of statistical models, especially for more complicated ones. Like Nimble, STAN compiles each model, but one big advantage is that it can store each compiled model in a `.Rds` file, which allows for repeated use without the need to re-compile each time. In R, STAN is used through the `rstan` package, which allows for data to be loaded through normal R data structures. The biggest practical difference in using STAN compared to other bayesian computational methods is in the model specification process. STAN requires a model object in the form of a `.stan` file, which is created in a very unique language that is more formal than in BUGS. A `.stan` file, either created through a character string in R or in a seperate text file, can be used by the `stan()`

function, which takes the file along with the data to perform the sampling. There is a very robust set of manuals for STAN users, compared to other methods.

## Greta

Greta, a very recent addition to the list of bayes computational methods, is made specifically for use in R, making the model and data specification process as easy as possible. It also uses Google TensorFlow to perform its calculations, and allows for more flexible use of CPUs and GPUs. Unfortunately, as of the writing of this, it is still in need of lots of development, and lacks an appropriate amount of documentation and support for users. Computationally, it is very slow compared to the other bayes methods, especially in simple cases. There are several types of semi-common models that greta seems unable to properly handle in model specification at the moment. Greta also attempts to make model specification very user-friendly, but uses some very common function names in the process that might override other often-used R functions. Greta does offer a model visualization tool that does not exist in any of the other Bayesian samplers, in that it can provide a graphical representation of the specified model through the use of nodes and arrows. This tool can be very useful in trying to explain how all the components of a model work together.

## Common Bayesian Computation Settings

### Chains

When using Markov Chain Monte Carlo methods to estimate Bayesian posteriors of parameters, the number of independent Markov chains to be run needs to be specified. It is recommended to use several chains to explore the space of the target distribution, as a single chain may not fully traverse the entire region on its own. For some simpler models that don't have irregular features in the target distribution, using one chain can suffice. Each sampler has an argument in its function call that specifies the number of chains. Some of the samplers have a default value if the number of chains isn't given.

### Iterations

The number of steps that each Markov chain will take is a very important setting to get right. If each chain is not set up to run long enough, the parameter space may not be fully explored and the estimates will be inaccurate. Checking convergence statistics can be helpful in determining if the number of iterations is sufficient. An argument specifying the number of iterations for each chain can be inputted in the function call for each sampler.

### Warmup

Since our inferences are based on the assumption that our draws are from the target distribution, it is wise to discard the samples at the beginning of each chain, since it takes time for the sampler to “settle down”. The terms “warmup” and “burn-in” are often used to indicate the number of iterations that we will discard at the beginning of each chain when calculating statistics about our posterior distribution. We can choose this number in the function call for each sampler. In **OpenBUGS**, **Nimble**, and **Jags**, the argument is referred to as `burnin`, but in **STAN** and **greta** it is called `warmup`. Some samplers (**greta**, **Jags**, **Nimble**) treat this warmup argument as a number of iterations to run before starting to track the sampling results, but others (**STAN**, **BUGS**) have this be the actual number of iterations to discard out of the total number of specified iterations, so it is important to know which designation a sampler uses.

## Thinning

In some cases, there is high correlation between sequential samples in a given chain, which can bias the resulting posterior inference, which can be diagnosed using some graphical procedures. In this situation, it is wise to “thin” the results by only looking at every  $i$ th observation from a chain, such as every 10th observation. While each sampler defaults to including every simulated posterior draw, a user can specify a “thinning rate”, which is the rate at which a sampler will keep memory of the observations. A thinning rate of 1 means that every observation will be kept, and a rate of 10 indicates that only every 10th observation will be used.

## Initial Value Specification

In order to start sampling, each Markov chain needs to have an initial value at which to start. These initial values can be either be specified, or can often be randomly generated based on the prior distributions for each parameter of interest. In many cases, it is easiest to have the initial values randomly generated by the sampler. However, in cases where the target distribution is very complex, it is possible that poorly chosen initial values can lead to a chain getting “stuck” in a difficult region, or diverging in the wrong direction. In these cases, a more strategic choice of initial values is advised. **Nimble** unfortunately requires for these values to be specified instead of randomly generated. Using R to randomly generate the values from the prior distributions can be used in this situation:

```
nimble_inits <- list(
  "alpha" = rnorm(1,0,1000),
  "beta" = rnorm(1,0,1000),
  "sigma" = runif(1,0,10)
)

nimbleMCMC("inits" = nimble_inits, ... )
```

If multiple chains are being used, a list containing a number of lists equalling number of chains can be used:

```
nimble_inits <- list(
  list(
    "alpha" = rnorm(1,0,1000),
    "beta" = rnorm(1,0,1000),
    "sigma" = runif(1,0,10)
  ),
  list(
    "alpha" = rnorm(1,0,1000),
    "beta" = rnorm(1,0,1000),
    "sigma" = runif(1,0,10)
  )
)
```

Alternatively, if a list containing the initial values for one chain is used, these values will be recycled for the remaining chains.

## Common Tricks for Model and Data Specification

When describing a bayesian model through code, there are instances where slightly different implementations are needed based on the sampler being used. Both in model and data specification, the user needs to understand the intricacies of the specific sampling tool being used in order to get the analysis to run properly.

This section provides some helpful tips for implementing models in each of the different samplers being examined.

## Model Specification

The most important distinction between each of the samplers is the way that each allows for model specification text to be inputted.

### OpenBUGS

In **OpenBUGS**, the earliest method created of the lot, the model specification code must be wrapped inside the expression `model{}`:

```
model{
  ...model code here...
}
```

Since we are using the **R2OpenBUGS** library to run **OpenBUGS** through an R session, a text file with the model code is passed to the `bugs()` function. We can specify the model in R and create the appropriate text file:

```
bugs_model <- function() {
  ...model code here...
}

bugs.file <- file.path(tempdir(), "model.txt")
write.model(bugs_model, bugs.file)

bugs("model" = bugs.file, ... )
```

### JAGS

The **JAGS** sampler was intended to allow for model specification code to be almost identical to **BUGS**, so the implementation is very similar. There are two common approaches, one using the **rjags** package, and one using the **runjags** package. We use the latter, which allows for the model code to be specified in a multi-line character string:

```
jags_model <- "
  model{
    ...model code here...
  }
"

run.jags("model" = jags_model, ... )
```

### Nimble

Model specification in **Nimble** is extremely similar to **BUGS** as well. Model code is wrapped inside the **nimble** package function `nimbleCode()`, and then fed into `nimbleMCMC()` sampler function:

```
nimble_model <- nimbleCode({
  ...model code here...
})

nimbleMCMC("code" = nimble_model, ... )
```

## STAN

The STAN model specification language is significantly different than other samplers, in that a .stan file is fed to the C++ compiler. A .stan file has a very unique model specification format, which divides code into “data”, “parameter”, and “model” blocks at a minimum. When needed, “transformed data” and “transformed parameters” blocks are also used. The basic format for model specification is provided below:

```
data {  
  ...data specified here...  
}  
  
parameters {  
  ...parameters specified here...  
}  
  
model {  
  ...model code here...  
}
```

Each data variable and parameter variable needs to be initialized in the appropriate block, along with associated domain limits. For example, if one were initializing a rate parameter for an exponential distribution called “lambda”, it would be initialized with a lower bound of 0:

```
parameters {  
  <lower=0> lambda;  
}
```

A character string specifying the path to the .stan file is then passed to `stan()`.

```
stan_file <- "file.stan"  
stan("file" = stan_file, ... )
```

## Data Specification and Input

Specifying data for the bayesian samplers in R is very straightforward. In OpenBUGS, JAGS, and STAN, a list containing each data object or constant can be passed to the function that calls the sampler. For example, in JAGS, the data might constructed like this:

```
n <- 10 # binomial n  
y <- rbinom(1, n, p)  
  
jags_data <- list(  
  "y" = y,  
  "n" = n  
)  
  
run.jags("data" = jags_data, ... )
```

In Nimble, the constants need to be specified in a separate list from the data.

```
n <- 10 # binomial n  
y <- rbinom(1, n, p)  
  
nimble_data <- list(  
  "y" = y  
)
```



```
nimble_constants <- list(
  "n" = n
)

nimbleMCMC("data" = nimble_data, "constants" = nimble_constants, ... )
```

## Monitoring Parameters

In OpenBUGS, JAGS, Nimble, and STAN, specifying parameters for the sampler to monitor is as easy as providing a character vector with the names of the parameters to monitor:

```
jags_monitors <- jags_monitor <- c("alpha", "beta")

run.jags("monitor" = jags_monitor, ... )
```

## Greta's Unique Specification Language

Implementing models in **greta** is very unique in that each data variable or parameter is initialized separately with its own line of R code. Each data object in R is redefined using the `as_data()` function in the **greta** package:

```
y <- as_data(rpois(5, theta))
```

Each parameter is defined using a function indicating the prior distribution and the chosen prior parameters.

```
theta <- gamma(3,1)
```

Once all the variables have been defined, the relationship between the data and the unknown parameters is established using the `distribution()` function:

```
distribution(y) <- poisson(theta)
```

Finally, the **greta** model is constructed using `model()` function, which takes each of the monitored parameters as arguments:

```
greta_model <- model(theta)

mcmc("model" = greta_model, ... )
```

## How to Specify Tricky Models

Each sampler has its own unique methods for using certain bayesian model details. We provide illustrations for some of the most common ones.

### Truncated Distribution

Using a truncated distribution is a common practice in bayesian modeling, especially when specifying priors for parameters with positive domains. The syntax for specifying a truncated distribution is slightly different for each sampling method. Below we are using a truncated normal distribution with a lower bound of zero and no upper bound:

```

## Truncated normal distribution in OpenBUGS
tau ~ dnorm(0,1) I(0, )

## Model Test for OpenBUGS using R2OpenBUGS
tau ~ dnorm(0,1) %_ I(0, )

## JAGS
tau ~ dnorm(0,1) I(0, )

## Nimble
tau ~ T(dnorm(0,1), 0,)

## Greta
tau <- normal(0,1, truncation = c(0,Inf))

## STAN
# Parameter block initialization
real<lower=0> tau;

# Model block specification
tau ~ normal(0,1);

# Alternative solution
tau ~ normal(0,1) T[0,]

```

## Flat Priors

Sometimes using a flat, diffuse prior is appropriate. Most of the samplers allow for improper flat priors, although JAGS does not. In STAN, a parameter has is given a flat, improper prior by default if it doesn't have a distribution assigned. Simply declaring the parameter will do the job.

```

## Flat prior in OpenBUGS
theta ~ dflat()

## Nimble
theta ~ dflat()

## Greta
theta <- variable()

## STAN
# Parameter block initialization
real theta;

```

## Distribution Parameter Order

Some samplers use different parameter orders for certain distributions.

Some use the  $\sigma$  parameterization of the normal distribution, where  $\sigma$  is the standard deviation, and others use  $\tau$  instead, where  $\tau = 1/\sigma^2$ . Among bayesian samplers, JAGS, OpenBUGS, and Nimble use the  $\tau$  parameterization, while STAN and greta use  $\sigma$ .

**Normal:** The  $N(\mu, \sigma)$  parameterization is used by STAN and greta, where  $\sigma$  is the standard deviation. The  $N(\mu, \tau)$  parameterization is used by JAGS, OpenBUGS, and Nimble, which uses the precision parameter

$\tau = 1/\sigma^2$ . The same parameterization guidelines are the same for the student-t distribution as well.

**Binomial:** The  $Bin(n, p)$  parameterization is used by **OpenBUGS**, **JAGS**, and **Nimble**, while the  $Bin(p, n)$  parameterization is used by **STAN** and **greta**.

**Gamma:** The  $Gam(\alpha, \lambda)$  parameterization, where  $\alpha$  is a shape parameter and  $\lambda$  is a rate parameter, is used by **OpenBUGS**, **JAGS**, **Nimble**, and **greta**. The  $Gam(\alpha, \beta)$  parameterization is used by **STAN**, using the scale parameter  $\beta = 1/\lambda$ .

**Weibull:** The  $Weibull(v, \lambda)$  parameterization is used by **OpenBUGS**, **JAGS**, **Nimble**, and **greta**, where  $\lambda$  is a shape parameter. The  $Weibull(v, \sigma)$  parameterization is used by **STAN**, where  $\sigma = 1/\lambda$ .

### Equations inside distribution parameters

In order to utilize calculated parameters inside **OpenBUGS**, a variable must be defined using the assignment operator that is the result of the calculation:

```
y_hat <- alpha + beta * x
y ~ dnorm(y_hat, tau)
```

Fortunately, the other samplers allow for variable calculations inside the actual parameter definition of a distribution function. Here is an example in **JAGS**:

```
y ~ dnorm(alpha + beta * x, tau)
```