

Hands-on training session 1

Introduction, revision and running basic models

Matt Denwood Giles Innocent

2020-02-12

Course Outline and Practicalities

Overview

Date/time:

- 19th February 2020
- 14.00 - 15.30

Teachers:

- Matt Denwood (presenter)
- Giles Innocent

Preparation

- Have we all looked at the preparation material?
- Do we all have R and JAGS installed?
- Brief demonstration of GitHub and where the files are. . .

Revision

Bayes Rule

Bayes' theorem is at the heart of Bayesian statistics. This states that:

$$P(\theta|Y) = \frac{P(\theta) \times P(Y|\theta)}{P(Y)}$$

Bayes Rule

Bayes' theorem is at the heart of Bayesian statistics. This states that:

$$P(\theta|Y) = \frac{P(\theta) \times P(Y|\theta)}{P(Y)}$$

Where: θ is our parameter value(s); Y is the data that we have observed; $P(\theta|Y)$ is the posterior probability of the parameter value(s) given the data and priors; $P(\theta)$ is the prior probability of the parameters BEFORE we had observed the data; $P(Y|\theta)$ is the likelihood of the data given the parameters value(s), as discussed above; and $P(Y)$ is the probability of the data, integrated over all parameter space.

Note that $P(Y)$ is rarely calculable except in the simplest of cases, but is a constant for a given model. So in practice we usually work with the following:

$$P(\theta|Y) \propto P(\theta) \times P(Y|\theta)$$

Note that $P(Y)$ is rarely calculable except in the simplest of cases, but is a constant for a given model. So in practice we usually work with the following:

$$P(\theta|Y) \propto P(\theta) \times P(Y|\theta)$$

Our Bayesian posterior is therefore always a combination of the likelihood of the data, and the parameter priors

But for more complex models the distinction between what is 'data' and 'parameters' can get blurred!

TODO

- Brief highlight of important points re convergence and effective sample size

Session 1a: Theory and application of MCMC

We have seen that we can write a Metropolis algorithm but this is complex and inefficient

There are a number of general purpose languages that allow us to define the problem in a fairly intuitive manner and leave the details to the software. Among the latter are BUGS (Bayesian inference Using Gibbs Sampling) JAGS (Just another Gibbs Sampler) and STAN (named in honour of Stanislaw Ulam, pioneer of the Monte Carlo method).

JAGS

JAGS is a declarative (non-procedural) programming language. That is the order of statements does not matter. When implemented the interpreter just considers the appropriate line of code to interpret in terms of likelihood and prior. That is why you can only define each variable (LHS) once.

Different ways to run JAGS from R:

- rjags
- runjags
- R2jags
- jagsUI
- See <http://runjags.sourceforge.net/quickjags.html> (also in the GitHub folder)

A simple JAGS model might look like this:

```
1  model{
2    # Likelihood part:
3    Positives ~ dbinom(prevalence, TotalTests)
4
5    # Prior part:
6    prevalence ~ dbeta(2, 2)
7
8    # Hooks for automatic integration with R:
9    #data# Positives, TotalTests
10   #monitor# prevalence
11   #inits# prevalence
12 }
```

There are two model statements:

```
1 Positives ~ dbinom(prevalence, TotalTests)
```

states that the number of *Positive* test samples is Binomially distributed with probability parameter *prevalence* and total trials *TotalTests*

```
1 prevalence ~ dbeta(2,2)
```

states that our prior probability distribution for the parameter *prevalence* is Beta(2,2)

These are very similar to the likelihood and prior functions defined in the preparatory exercise

The other lines in this model are automated hooks to allow the `runjags` package to make sure that the correct R objects are passed to R as data, and that we get samples for the parameter we want to monitor. These statements are ignored by JAGS.

One of the advantages of using something like JAGS or STAN is that much of the detail of the coding is shielded from us which means that the code is accessible to a wider variety of audiences. Running this model is much easier (and more efficient) than the Metropolis algorithm code.

To run this model, copy/paste the code above into a new text file called “basicjags.bug” in the same folder as your current working directory.

Then run:

```
1 library('runjags')

1 ## Attaching runjags (version 2.0.4-6) and setting user-specified
  ↪ options

1 # data to be retrieved by runjags:
2 Positives <- 7
3 TotalTests <- 10
4
5 # initial values to be retrieved by runjags:
6 prevalence <- list(chain1=0.05, chain2=0.95)
```

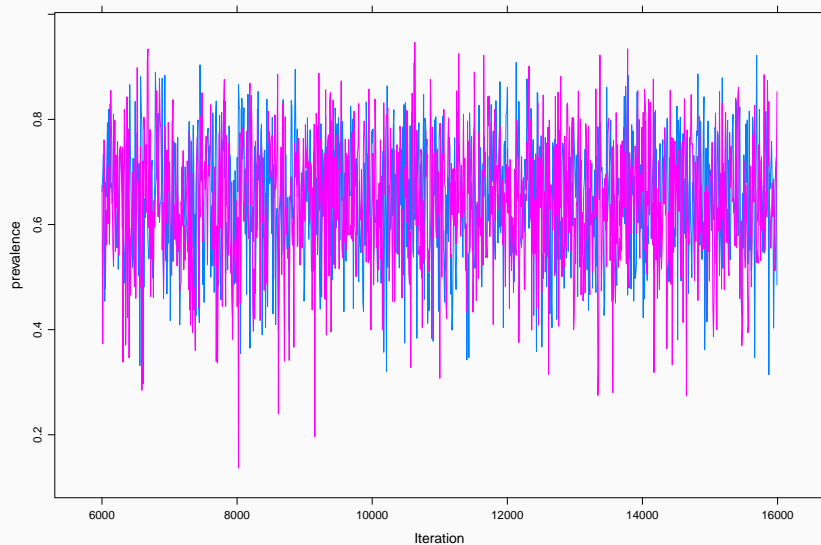
```
1 results <- run.jags('basicjags.bug', n.chains=2, burnin=5000,  
  ↪ sample=10000)  
  
1 ## Loading required namespace: rjags  
  
1 ## Compiling rjags model...  
2 ## Calling the simulation using the rjags method...  
3 ## Adapting the model for 1000 iterations...  
4 ## Burning in the model for 5000 iterations...  
5 ## Running the model for 10000 iterations...  
6 ## Simulation complete  
7 ## Calculating summary statistics...  
8 ## Calculating the Gelman-Rubin statistic for 1  
9 ## variables....  
10 ## Finished running the simulation
```

You should see some updates then it will be finished

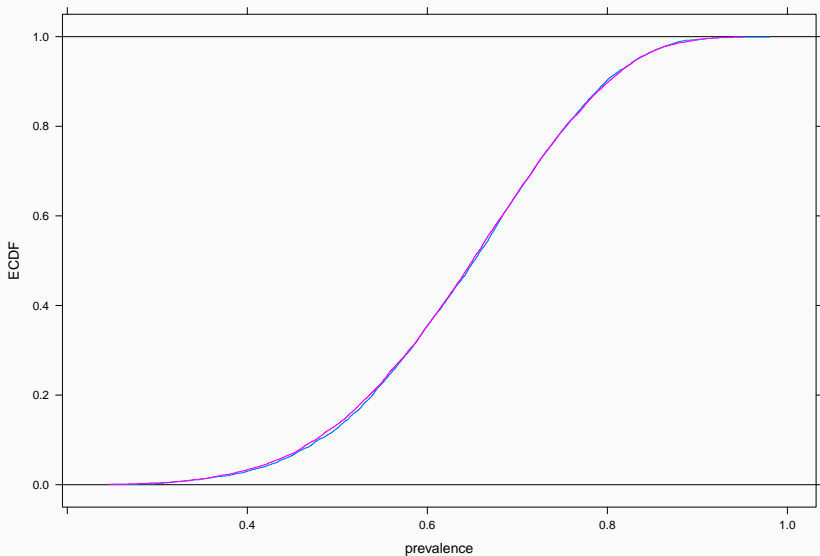
First check the plots for convergence:

```
1 plot(results)
```

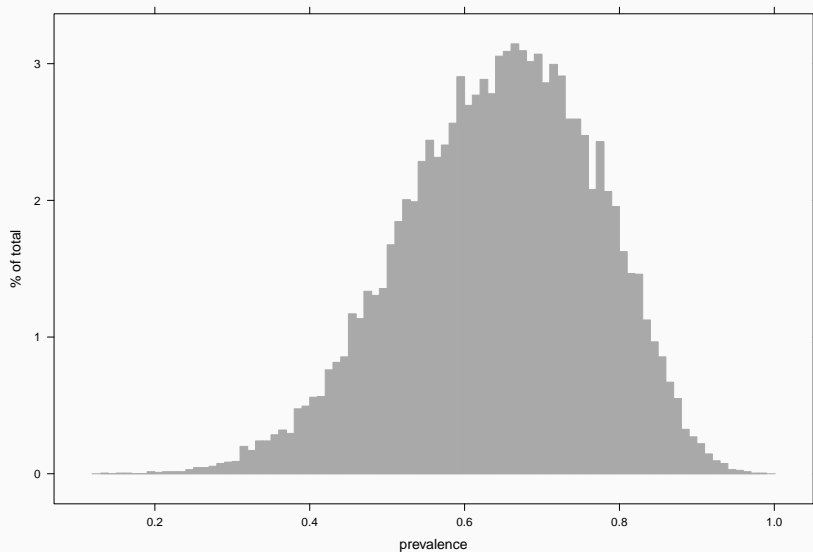
Trace plots: the two chains should be stationary:



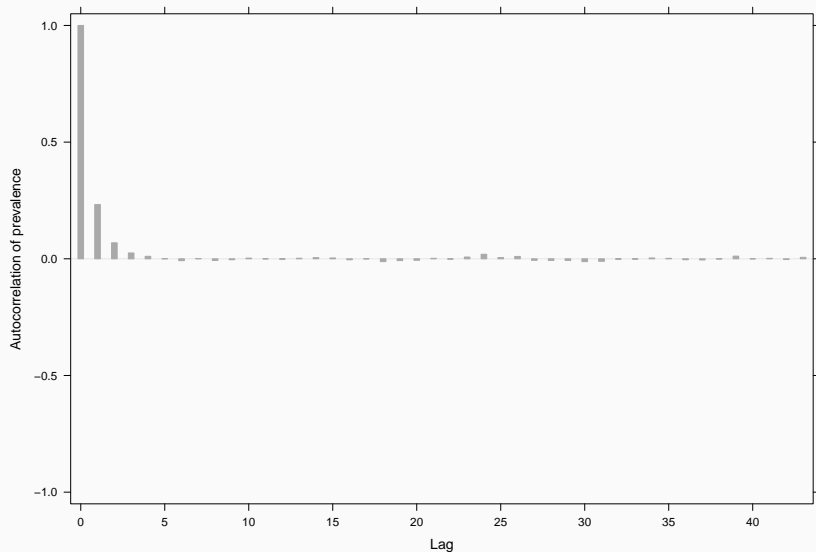
ECDF plots: the two chains should be very close to each other:



Histogram of the combined chains should appear smooth:



Autocorrelation plot tells you how well behaved the model is:



Then check the effective sample size (S_{Seff}) and Gelman-Rubin statistic (psrf):

```
1 results

1 ##
2 ## JAGS model summary statistics from 20000 samples (chains = 2;
   ↳ adapt+burnin = 6000):
3 ##
4 ##           Lower95  Median Upper95      Mean      SD      Mode
5 ## prevalence 0.40369 0.65081 0.87298 0.64327 0.12392 0.66933
6 ##
7 ##           MCerr MC%ofSD SSeff      AC.10  psrf
8 ## prevalence 0.001139      0.9 11837 0.0028987 1.0001
9 ##
10 ## Total time taken: 0.6 seconds
```

Reminder: we want $psrf < 1.05$ and $S_{Seff} > 1000$. If these both look OK then you can use the posterior summary statistics.

Exercise

- Run this model yourself in JAGS
- Change the initial values for the two chains and make sure it doesn't affect the results
- Reduce the burnin length - does this make a difference?
- Change the sample length - does this make a difference?

Optional Exercise

- Change the number of chains to 1 and 4
 - Remember that you will also need to change the initial values
 - What affect does having different numbers of chains have?
 - Try using the `run.jags` argument `method='parallel'` - what affect does this have?

Session 1b: Working with basic models (apparent prevalence)

Other runjags options

There are a large number of other options to runjags. Some highlights:

- The method can be parallel or background or bgparallel
- You can use `extend.jags` to continue running an existing model (e.g. to increase the sample size)
- You can use `coda::as.mcmc.list` to extract the underlying MCMC chains
- Use the `summary()` method to extract summary statistics
 - See `'?summary.runjagsand?runjagsclass'` for more information

Using embedded character strings

- For simple models we might not want to bother with an external text file. Then we can do:

```
1  mt <- "  
2  model{  
3    Positives ~ dbinom(prevalence, TotalTests)  
4    prevalence ~ dbeta(2, 2)  
5  
6    #data# Positives, TotalTests  
7    #monitor# prevalence  
8    #inits# prevalence  
9  }  
10 "  
11  
12 results <- run.jags(mt, n.chains=2)  
  
1  ## Finished running the simulation
```

- But this makes the model harder to read, so I would generally

Setting the RNG seed

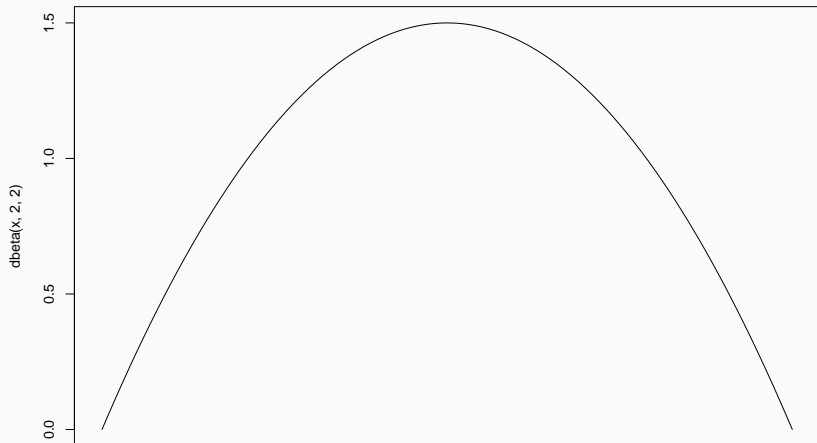
- If we want to get numerically replicable results we need to add `.RNG.name` and `.RNG.seed` to the initial values, and an additional `#modules#` `lecuyer` directive:

```
1  mt <- "  
2  model{  
3    Positives ~ dbinom(prevalence, TotalTests)  
4    prevalence ~ dbeta(2, 2)  
5  
6    #data# Positives, TotalTests  
7    #monitor# prevalence  
8    #inits# prevalence, .RNG.name, .RNG.seed  
9    #modules# lecuyer  
10 }  
11 "  
12 .RNG.name <- "lecuyer::RngStream"  
13 .RNG.seed <- list(chain1=1, chain2=2)  
14 results <- run.jags(mt, n.chains=2)
```

A different prior

- A quick way to see what the distribution of a prior looks like is via e.g.:

```
1 curve(dbeta(x, 2, 2), from=0, to=1)
```



- Let's change the prior we are using to `dbeta(1,1)`:

```
1  model{
2    Positives ~ dbinom(prevalence, TotalTests)
3    prevalence ~ dbeta(1, 1)
4
5    # Hooks for automatic integration with R:
6    #data# Status, TotalTests
7    #monitor# prevalence
8    #inits# prevalence
9  }
```


An Equivalent Model

- We could equivalently specify the model as:

```
1  model{
2    # Likelihood part:
3    for(i in 1:TotalTests){
4      Status[i] ~ dbern(prevalence)
5    }
6
7    # Prior part:
8    prevalence ~ dbeta(1, 1)
9
10   # Hooks for automatic integration with R:
11   #data# Status, TotalTests
12   #monitor# prevalence
13   #inits# prevalence
14 }
```

- But we need the data in a different format: a vector of 0/1 rather than total positives!

A GLM Model

- Another option is:

```
1  model{
2    # Likelihood part:
3    for(i in 1:TotalTests){
4      Status[i] ~ dbern(predicted[i])
5      logit(predicted[i]) <- intercept
6    }
7
8    # Prior part:
9    intercept ~ dnorm(0, 10^-6)
10
11   # Derived parameter:
12   prevalence <- ilogit(intercept)
13
14   # Hooks for automatic integration with R:
15   #data# Status, TotalTests
16   #monitor# intercept, prevalence
17   #inits# intercept
18 }
```

- This is the start of a generalised linear model, where we could add covariates at individual animal level.
- We introduce a new distribution `dnorm()` - notice this is mean and precision, not mean and sd!
- For a complete list of the distributions available see:
 - <https://sourceforge.net/projects/mcmc-jags/files/Manuals/4.x/>
 - This document is also provided on the GitHub repository
- However, notice that the prior is specified differently. . .

Exercise

- Run the original version of the model and the GLM version of the model and compare results with the same data
- Now try a larger sample size: e.g. 70 positives out of 100 tests
 - are the posteriors from the two models more or less similar than before?
- Now try running the GLM model with a prior of `dnorm(0, 0.33)` (and the original data) - does this make a difference?

Optional Exercise

Another way of comparing different priors is to run different models with no data - as there is no influence of a likelihood, the posterior will then be identical to the priors (and the model will run faster).

One way to do this is to make all of the response data (i.e. either Positives or Status) missing. Try doing this for the following three models, and compare the priors for prevalence:

- The original model with prior prevalence $\sim \text{dbeta}(1,1)$
- The GLM model with prior intercept $\sim \text{dnorm}(0, 10^{-6})$
- The GLM model with prior intercept $\sim \text{dnorm}(0, 0.33)$

Session 1c: Basics of latent-class models (imperfect test)

Imperfect tests

- Up to now we have ignored issues of diagnostic test sensitivity and specificity
- Usually, however, we do not have a perfect test, so we do not know how many are truly positive or truly negative, rather than just testing positive or negative.
- But we know that:

$$Prev_{obs} = (Prev_{true} \times Se) + ((1 - Prev_{true}) \times (1 - Sp))$$

$$\implies Prev_{true} = \frac{Prev_{obs} - (1 - Sp)}{Se - (1 - Sp)}$$

Model Specification

- We can write this model as:

```
1  model{
2    Positives ~ dbinom(obsprev, TotalTests)
3    obsprev <- (prevalence * se) + ((1-prevalence) * (1-sp))
4
5    prevalence ~ dbeta(1, 1)
6    se ~ dbeta(1, 1)
7    sp ~ dbeta(1, 1)
8
9    #data# Positives, TotalTests
10   #monitor# prevalence, obsprev, se, sp
11   #inits# prevalence, se, sp
12 }
```


- And run it:

```
1 prevalence <- list(chain1=0.05, chain2=0.95)
2 se <- list(chain1=0.5, chain2=0.99)
3 sp <- list(chain1=0.5, chain2=0.99)
4 Positives <- 70
5 TotalTests <- 100
6 results <- run.jags('basicimperfect.bug', n.chains=2, burnin=0,
  ↪ sample=10000)

1 ## Finished running the simulation
```

[Remember to check convergence and effective sample size!]

	Lower95	Median	Upper95	SSeff	psrf
prevalence	0.010	0.485	0.944	2130	1.000
obsprev	0.604	0.694	0.778	21339	1.000
se	0.141	0.688	1.000	1434	1.002
sp	0.000	0.306	0.844	1444	1.000

- What do these results tell us?

	Lower95	Median	Upper95	SSeff	psrf
prevalence	0.010	0.485	0.944	2130	1.000
obsprev	0.604	0.694	0.778	21339	1.000
se	0.141	0.688	1.000	1434	1.002
sp	0.000	0.306	0.844	1444	1.000

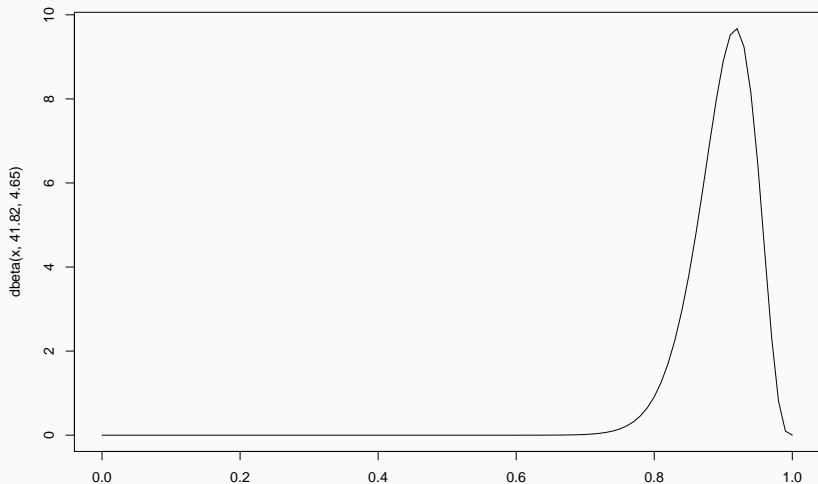
- What do these results tell us?
- We can estimate the observed prevalence quite well
- But not the prevalence, se or sp! The model is unidentifiable.

Priors

- We cannot estimate se, sp and prevalence simultaneously
- We need strong priors for se and sp
- We can use the PriorGen package to generate Beta priors based on published results, for example:

```
1 PriorGen::findbeta(themean=0.9, percentile = 0.975,  
  ↪ percentile.value = 0.8)  
  
1 ## [1] "The desired Beta distribution that satisfies the  
  ↪ specified conditions is: Beta( 41.82 4.65 )"  
2 ## [1] "Here is a plot of the specified distribution."  
3 ## [1] "Descriptive statistics for this distribution are:"  
4 ##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.  
5 ## 0.6556 0.8751 0.9057 0.9004 0.9320 0.9939  
6 ## [1] "Verification: The percentile value 0.8 corresponds to the  
  ↪ 0.025 th percentile"
```

```
1 qbeta(c(0.025, 0.5, 0.975), 41.82, 4.65)
1 ## [1] 0.7999144 0.9056630 0.9677283
1 curve(dbeta(x, 41.82, 4.65), from=0, to=1)
```



Exercise

- Find beta distribution priors for sensitivity = 0.9 (95% CI: 0.85 - 0.95) and specificity = 0.95 (95%CI: 0.92-0.97)
- Look at these distributions using curve and qbeta
- Modify the imperfect test model using these priors and estimate prevalence now

Optional Exercise

Run the same model with se and sp fixed to the mean estimate - how does this affect CI for prevalence?

Run the same model with se and sp fixed to 1 - how does this affect estimates and CI for prevalence?

Solution

```
1 PriorGen::findbeta(themean=0.9, percentile = 0.975,  
  ↪ percentile.value = 0.85)  
  
1 ## [1] "The desired Beta distribution that satisfies the  
  ↪ specified conditions is: Beta( 148.43 16.49 )"  
2 ## [1] "Here is a plot of the specified distribution."  
3 ## [1] "Descriptive statistics for this distribution are:"  
4 ##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.  
5 ##  0.7942  0.8852  0.9018  0.9002  0.9169  0.9613  
6 ## [1] "Verification: The percentile value 0.85 corresponds to  
  ↪ the 0.025 th percentile"  
  
1 qbeta(c(0.025, 0.5, 0.975), 148.43, 16.49)  
  
1 ## [1] 0.8500135 0.9016287 0.9408480  
  
1 curve(dbeta(x, 148.43, 16.49), from=0, to=1)
```


Summary