

# Group Work 05

```
library(bayesrules) # R package for our textbook
library(tidyverse) # Collection of packages for tidying and plotting data
library(janitor) # Helper functions like tidy and tabyl
library(rstan)
library(broom)
```

## i Note

The Monday problems are a selection from BR Exercises 8.15-8.17, packaged in a different format

## i Note

Make sure to remove `#| eval: false` from any code chunks if you copy the whole document using the “code” button

## 1 An exact analysis

This problem uses the `pulse_of_the_nation` data from slides13

```
data("pulse_of_the_nation")
pulse_of_the_nation %>%
  count(climate_change)
```

```
# A tibble: 3 x 2
  climate_change      n
  <fct>                <int>
1 Not Real At All        150
2 Real and Caused by People    655
3 Real but not Caused by People   195
```

- (a) Bayesians think of the entire Beta(151,852) posterior pdf as an estimate for  $\pi$ . But for communication purposes, it can be useful to report what values of  $\pi$  are typical. Identify and calculate two possible posterior point estimates.
- (b) The posterior estimates above merely capture the typical posterior  $\pi$  value, thus miss the bigger picture. It's important to supplement these estimates with a posterior credible interval. (Bayesians use "credible" instead of "confidence") Calculate a 95% posterior credible interval for  $\pi$ . Revisiting a plot of the posterior might spark some ideas, and you'll need some R code of the `rXXX`, `qXXX`, `dXXX`, `pXXX` variety.
- (c) How can we interpret this interval, (a,b)?
- (d) What would a strict frequentist say if you asked them "What's the probability that  $\pi$  lies within the Bayesian credible interval?"
- (e) A researcher claims that more than 13% of people don't believe in climate change. Using your interval from (b), what do you think about this claim?
- (f) Calculate and interpret a posterior probability that helps you test this claim. You'll need some R code of the `rXXX`, `qXXX`, `dXXX`, `pXXX` variety.
- (g) There's no common cut-off / threshold (eg: 0.05) for interpreting Bayesian posterior probabilities, hence no binary conclusion. Better yet, Bayesian conclusions are more holistic and nuanced. With this in mind, summarize your conclusions about our hypothesis.

## 2 Bayes Factor

### i Posterior Odds

The **posterior odds** for a hypothesis test  $H_0$  against  $H_a$  after observing data  $Y = y$  is

$$\text{posterior odds} = \frac{P(H_a|Y = y)}{P(H_0|Y = y)}$$

### i Posterior Odds

The **prior odds** for a hypothesis test  $H_0$  against  $H_a$  is

$$\text{posterior odds} = \frac{P(H_a)}{P(H_0)}$$

## i Bayes Factor

In a hypothesis test of two competing hypotheses,  $H_a$  vs  $H_0$ , the **Bayes Factor** is an odds ratio for  $H_a$ :

$$\text{Bayes Factor} = \frac{\text{posterior odds}}{\text{prior odds}} = \frac{P(H_a|Y)/P(H_0|Y)}{P(H_a)/P(H_0)}$$

Calculate the **prior odds**, **posterior odds**, and **Bayes Factor** for the researcher's claim that "more than 13% of people don't believe in climate change". What do these numbers tell you?

## 3 BR Exercise 8.21

Now, let's explore how we can do estimation and hypothesis testing with an approximate posterior sample from MCMC.

- Load the following packages, and run the following R/stan code to fit an MCMC approximation for the same problem. Make sure you understand all of the pieces of the model fitting code.

```
library(rstan)
library(bayesrules)
library(bayesplot)
library(broom)
```

```
# Define the Beta-Binomial model in rstan notation
climate_model <- "
  data {
    real<lower=0> alpha;
    real<lower=0> beta;
    int<lower=1> n;
    int<lower=0, upper=n> Y;
  }

  parameters {
    real<lower=0, upper=1> pi;
  }

  model {
    Y ~ binomial(n, pi);
    pi ~ beta(alpha, beta);
  }
```

```

"
# Set the random number seed
set.seed(84735)

# SIMULATE the posterior
climate_sim <- stan(
  model_code = climate_model,
  data = list(alpha = 1, beta = 2, Y = 150, n = 1000),
  chains = 4, iter = 5000*2)

```

- (b) Give the following plots a quick peak. Do you see any red flags?

```

mcmc_trace(climate_sim, pars = "pi")
mcmc_dens_overlay(climate_sim, pars = "pi")
mcmc_dens(climate_sim, pars = "pi")

```

- (c) The four chains in `climate_sim` are currently stored as an array. Use the code below to store all the chains in a single data frame.

```

# Store the array of 4 chains in 1 data frame
climate_chains <- as.data.frame(
  climate_sim,
  pars = "lp__", include = FALSE)

# Check out the results
dim(climate_chains)
head(climate_chains)

```

- (d) Recall your exact posterior point estimates of  $\pi$  from earlier. Estimate these posterior features using your MCMC simulation. (How accurate are these estimates?) NOTE: The `sample_mode()` function in `{bayesrules}` calculates the mode of a sample.

```

climate_chains %>%
  summarize(___)

```

- (e) Recall your exact 95% posterior credible interval for  $\pi$  from earlier. Estimate this interval using your MCMC simulation. (How accurate is this estimate?)
- (f) Recall your exact analysis of the claim that more than 13% of people don't believe in climate change, i.e.  $\pi > 0.13$ . Estimate the posterior probability of this claim using your MCMC simulation. (How accurate is this estimate?)
- (g) Play around with the following shortcut functions that can address some, but not all, of our posterior questions. These are applied directly to `climate_sim`, not `climate_chains`. Take notes what they do (leaving comments in your .qmd would be sufficient!)

```
# What is the estimate? The posterior mean, median, or mode?
tidy(climate_sim, conf.int = TRUE, conf.level = 0.95)

#
mcmc_areas(climate_sim, pars = "pi", prob = 0.95)
```

## 4 BR Exercise 8.18 (slightly modified)

We'll continue with the `pulse_of_the_nation` beta-binomial example from last time.

- (a) Suppose we plan to survey 20 more people and want to predict  $Y'$ , the number who don't believe in climate change. Use your posterior point estimates from Q1 to provide point estimates for  $Y'$
- (b) Use your MCMC simulation (`climate_chains`) to approximate the posterior predictive model of  $Y'$ , the number (out of  $n = 20$ ) that don't believe in climate change.

 Tip

- The 20,000 chain values capture one source of variability: the posterior variability in  $\pi$ . To capture the other source, think about the variability in  $Y'|\pi$ .
- Since this is a random process, make sure to add `set.seed` at the beginning of your code chunk

- (c) Construct a graphical summary of the (approximated) posterior predictive model of  $Y'$ . Remember that  $Y'$  should be integer values from 0 to 20. If it's not, go back to the previous part and try another technique.
- (d) Approximate an 80% prediction interval for  $Y'$
- (e) Approximate the probability that more than 5 of the 20 people won't believe in climate change

## 5 BR Exercise 8.11 (a) - (c)

As discussed in Section 8.3, it is sometimes possible to derive an exact posterior predictive model. Such is the case with the conjugate models we have studied thus far. To begin, suppose we observe  $Y = y$  successes in  $n$  trials where  $Y|\pi \sim Bin(n, \pi)$  and  $\pi \sim Beta(\alpha, \beta)$ .

- (a) Write down the posterior pdf of  $\pi|Y$  (should depend on  $y, n, \alpha, \beta, \pi$ )
- (b) Suppose we conduct  $n'$  new trials (where  $n'$  might differ from our original number of trials  $n$ ) and let  $Y' = y'$  be the observed number of successes in these new trials. Identify the conditional pmf of  $Y'|\pi$ ,  $f(Y'|\pi)$ . (should depend on  $y', \pi, n'$ )

(c) Prove that the posterior predictive distribution of  $Y'$ ,  $f(y'|y)$ , is given by:

$$f(y'|y) = \int f(y'|\pi)f(\pi|y)d\pi = \binom{n'}{y'} \frac{\Gamma(\alpha + \beta + n)}{\Gamma(\alpha + y)\Gamma(\beta + n - y)} \frac{\Gamma(\alpha + y + y')\Gamma(\beta + n - y + n' - y')}{\Gamma(\alpha + \beta + n + n')}$$