

3 - Sampling the Imaginary

Antonio Jurlina

07/06/2020

Easy

These problems use the samples from the posterior distribution for the globe tossing example. This code will give you a specific set of samples, so that you can check your answers exactly.

```
n <- 1000
n_success <- 6
n_trials <- 9
n_samples <- 1e4

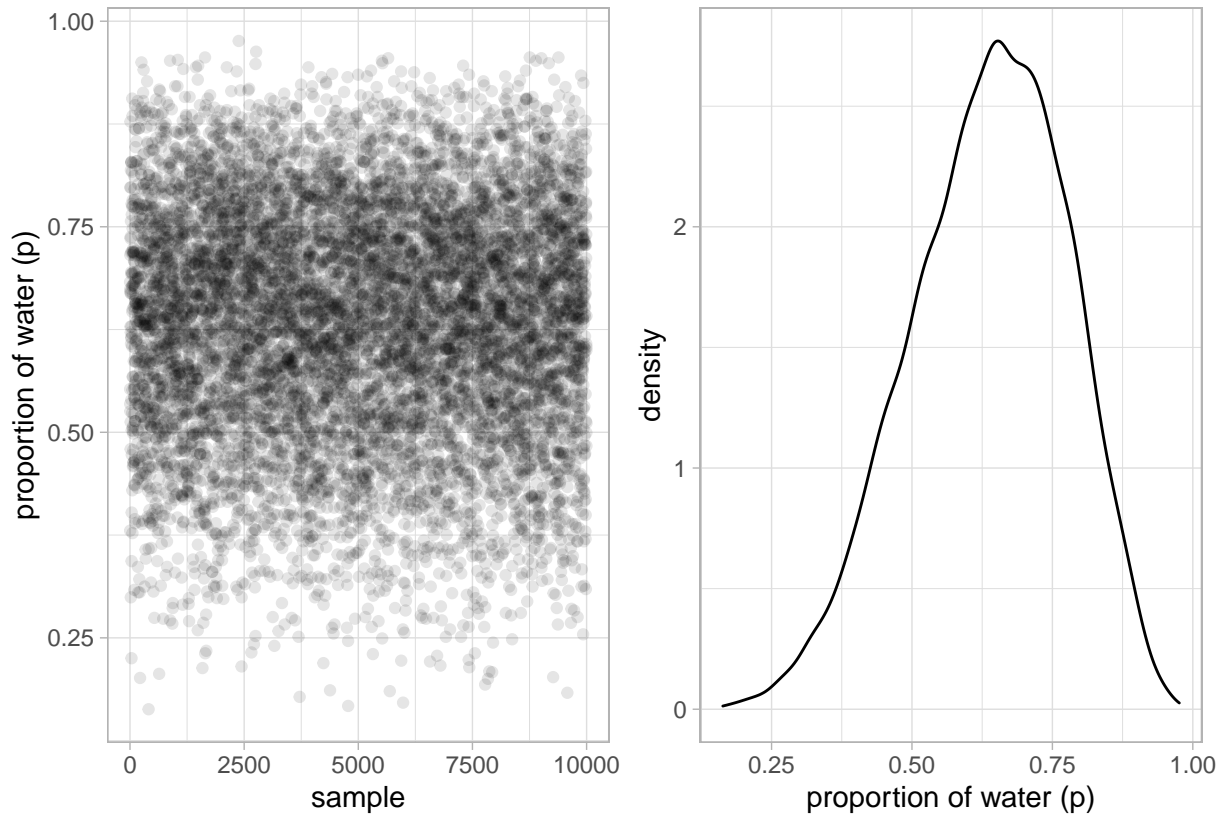
data <- tibble(p_grid = seq(from = 0, to = 1, length.out = n),
               prior = 1) %>%
  mutate(likelihood = dbinom(n_success, size = n_trials, prob = p_grid),
         posterior = (likelihood * prior) / sum(likelihood * prior))

samples <- data %>%
  sample_n(size = n_samples, weight = posterior, replace = TRUE)

plot1 <- samples %>%
  mutate(sample = 1:n()) %>%
  ggplot(aes(x = sample, y = p_grid)) +
  geom_point(alpha = 0.1) +
  theme_light() +
  ylab("proportion of water (p)")

plot2 <- samples %>%
  ggplot(aes(x = p_grid)) +
  geom_density() +
  theme_light() +
  xlab("proportion of water (p)")

plot1 + plot2
```



Use the values in samples to answer the questions that follow.

3E1. How much posterior probability lies below $p = 0.2$?

```
samples %>%
  summarise(sum = mean(p_grid < .2) * 100) %>%
  pull() %>%
  cat("%", sep = "")
```

```
## 0.08%
```

3E2. How much posterior probability lies above $p = 0.8$?

```
samples %>%
  summarise(sum = mean(p_grid > .8) * 100) %>%
  pull() %>%
  cat("%", sep = "")
```

```
## 11.49%
```

3E3. How much posterior probability lies between $p = 0.2$ and $p = 0.8$?

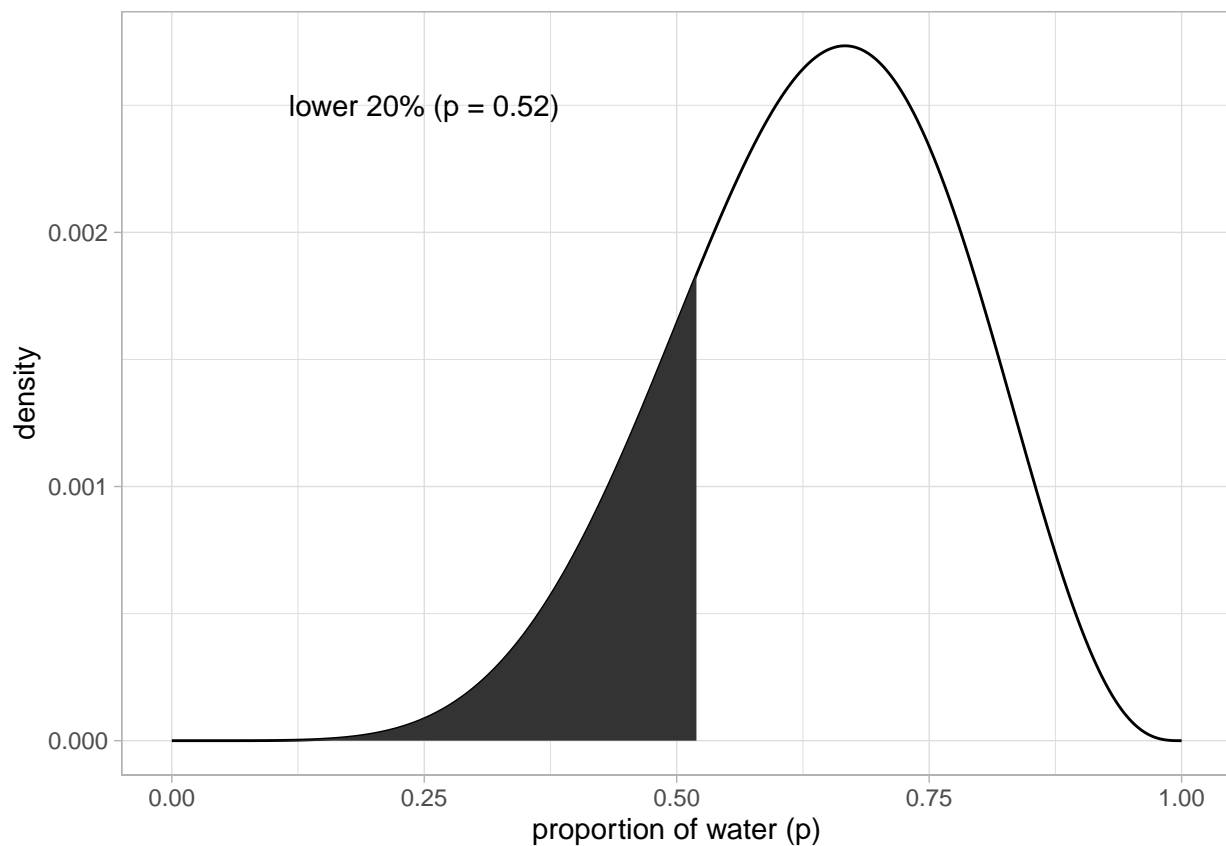
```
samples %>%
  summarise(sum = mean(p_grid > 0.2 & p_grid < .8) * 100) %>%
  pull() %>%
  cat("%", sep = "")
```

```
## 88.43%
```

3E4. 20% of the posterior probability lies below which value of p ?

```
q <- quantile(samples$p_grid, 0.2) %>%
  round(2)

data %>%
  ggplot(aes(x = p_grid)) +
  geom_line(aes(y = posterior)) +
  geom_ribbon(data = data %>% filter(p_grid < q),
            aes(ymin = 0, ymax = posterior)) +
  annotate(geom = "text",
          x = 0.25, y = 0.0025,
          label = paste0("lower 20% (p = ", q, ")")) +
  labs(x = "proportion of water (p)",
       y = "density") +
  theme_light()
```



3E5. 20% of the posterior probability lies above which value of p?

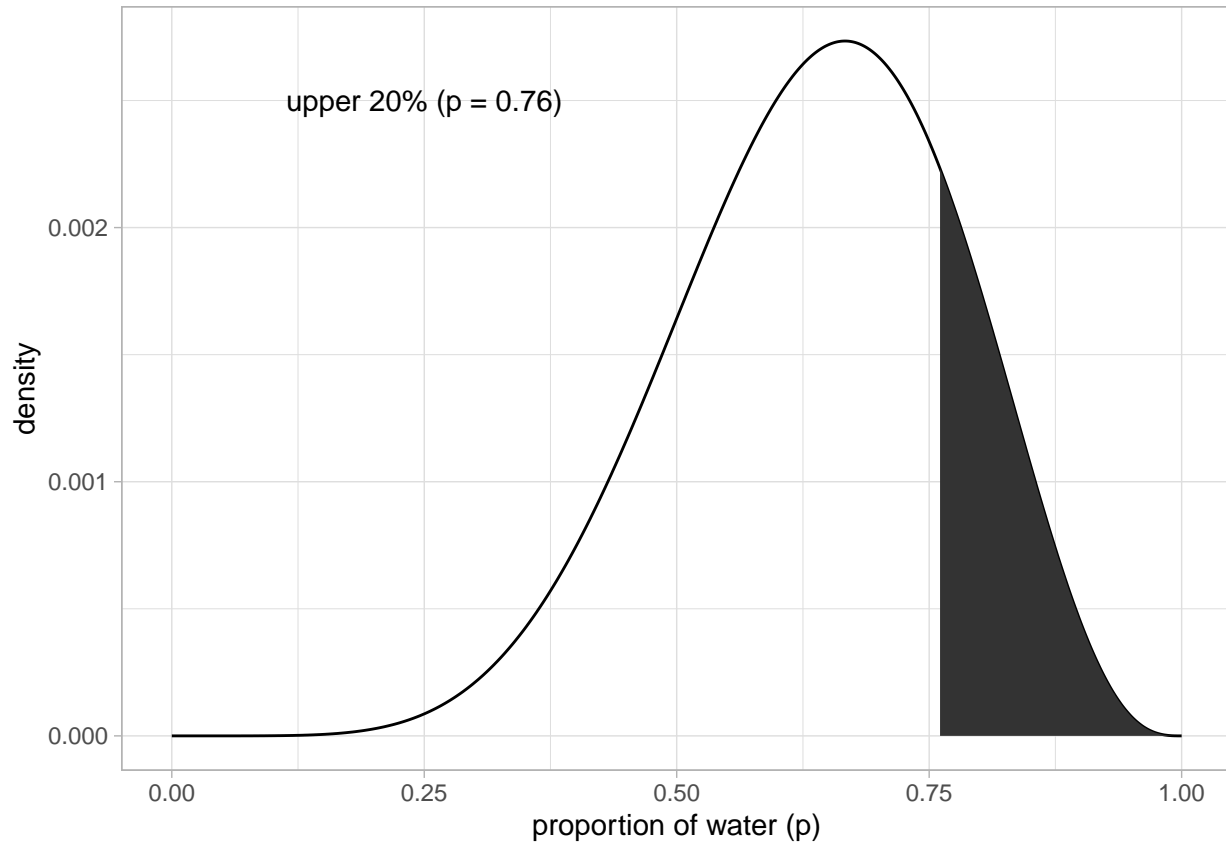
```
q <- quantile(samples$p_grid, 0.8) %>%
  round(2)

data %>%
  ggplot(aes(x = p_grid)) +
  geom_line(aes(y = posterior)) +
  geom_ribbon(data = data %>% filter(p_grid > q),
            aes(ymin = 0, ymax = posterior)) +
  annotate(geom = "text",
          x = 0.25, y = 0.0025,
```

```

    label = paste0("upper 20% (p = ", q, ")") +
  labs(x = "proportion of water (p)",
       y = "density") +
  theme_light()

```



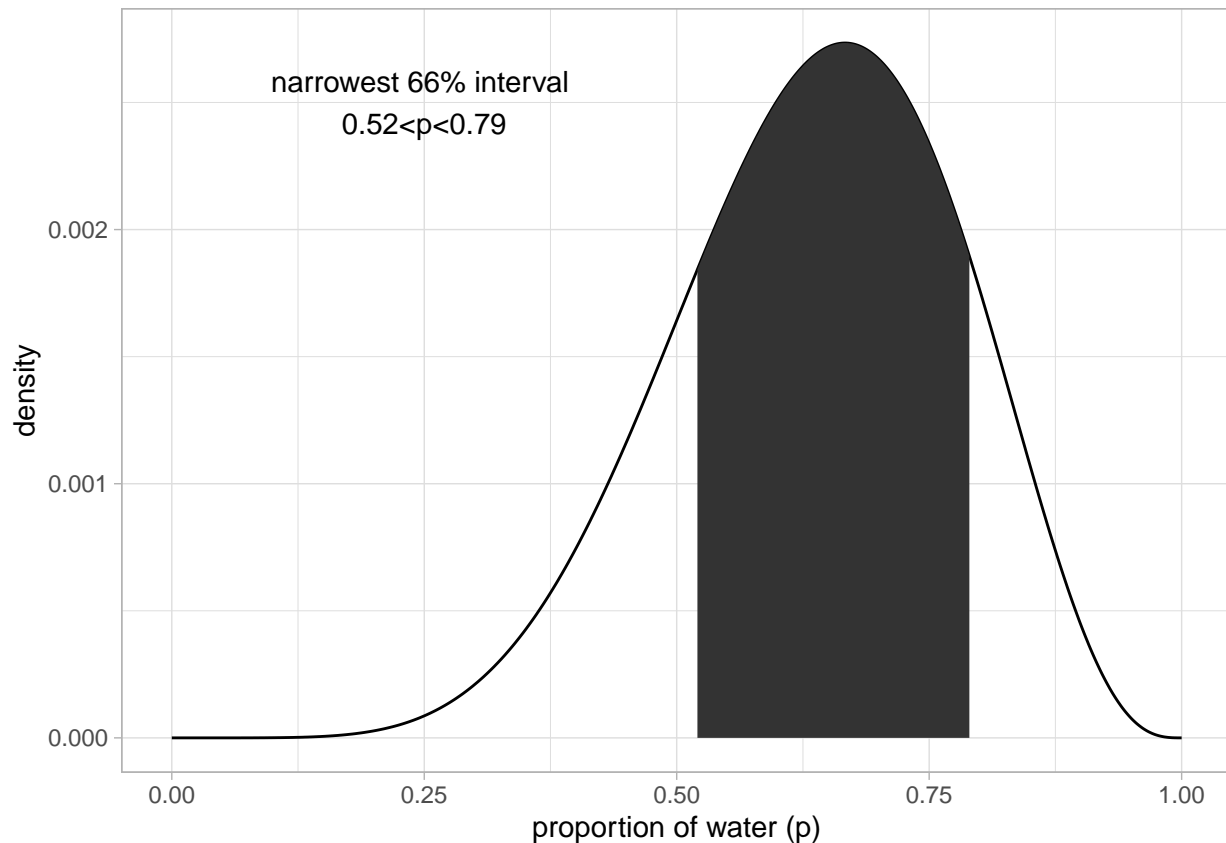
3E6. Which values of p contain the narrowest interval equal to 66% of the posterior probability?

```

HPDI(samples$p_grid, prob = 0.66) %>%
  round(2) %>% c(lower, upper)

data %>%
  ggplot(aes(x = p_grid)) +
  geom_line(aes(y = posterior)) +
  geom_ribbon(data = data %>% filter(p_grid > lower, p_grid < upper),
            aes(ymin = 0, ymax = posterior)) +
  annotate(geom = "text",
          x = 0.25, y = 0.0025,
          label = paste0("narrowest 66% interval \n",
                        lower, "<p<", upper)) +
  labs(x = "proportion of water (p)",
       y = "density") +
  theme_light()

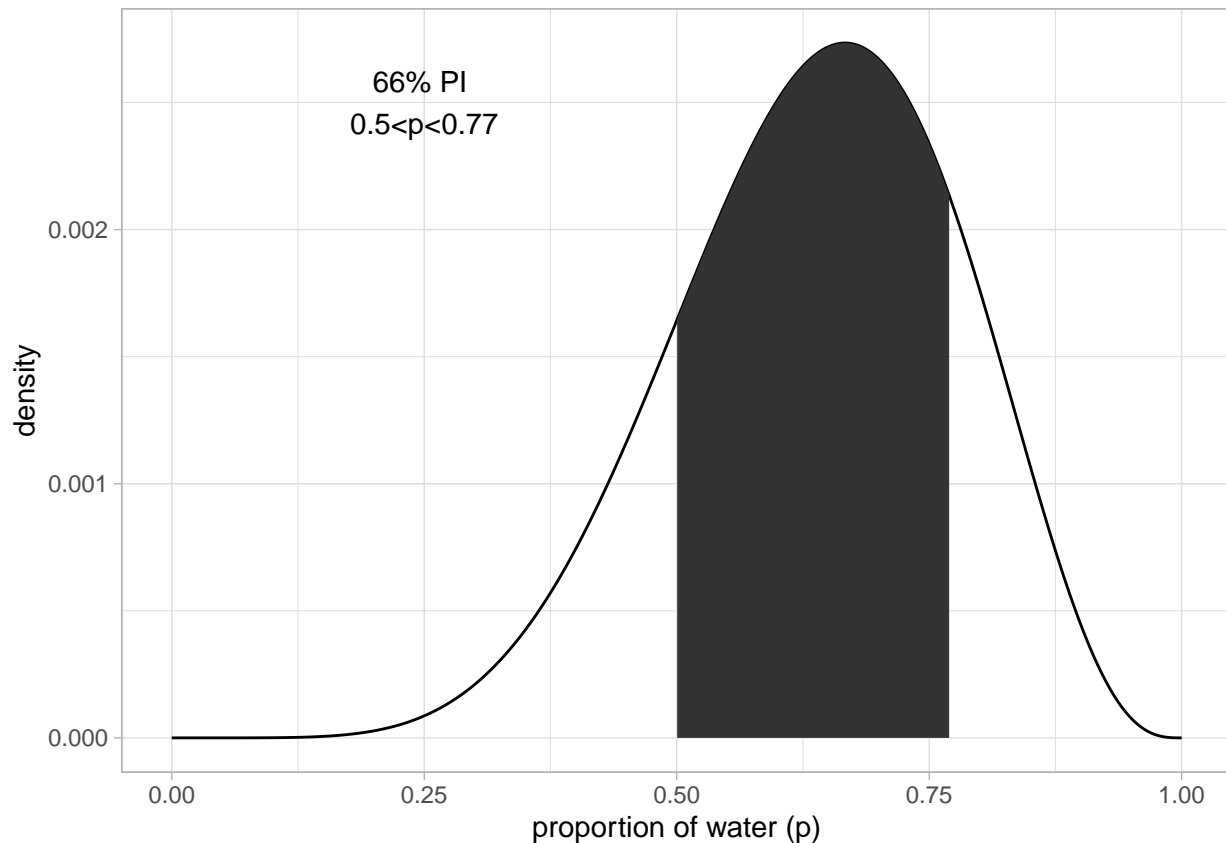
```



3E7. Which values of p contain 66% of the posterior probability, assuming equal posterior probability both below and above the interval?

```
PI(samples$p_grid, prob = 0.66) %>%
  round(2) %>% c(lower, upper)

data %>%
  ggplot(aes(x = p_grid)) +
  geom_line(aes(y = posterior)) +
  geom_ribbon(data = data %>% filter(p_grid > lower, p_grid < upper),
            aes(ymin = 0, ymax = posterior)) +
  annotate(geom = "text",
          x = 0.25, y = 0.0025,
          label = paste0("66% PI \n",
                        lower, "<p<", upper)) +
  labs(x = "proportion of water (p)",
       y = "density") +
  theme_light()
```



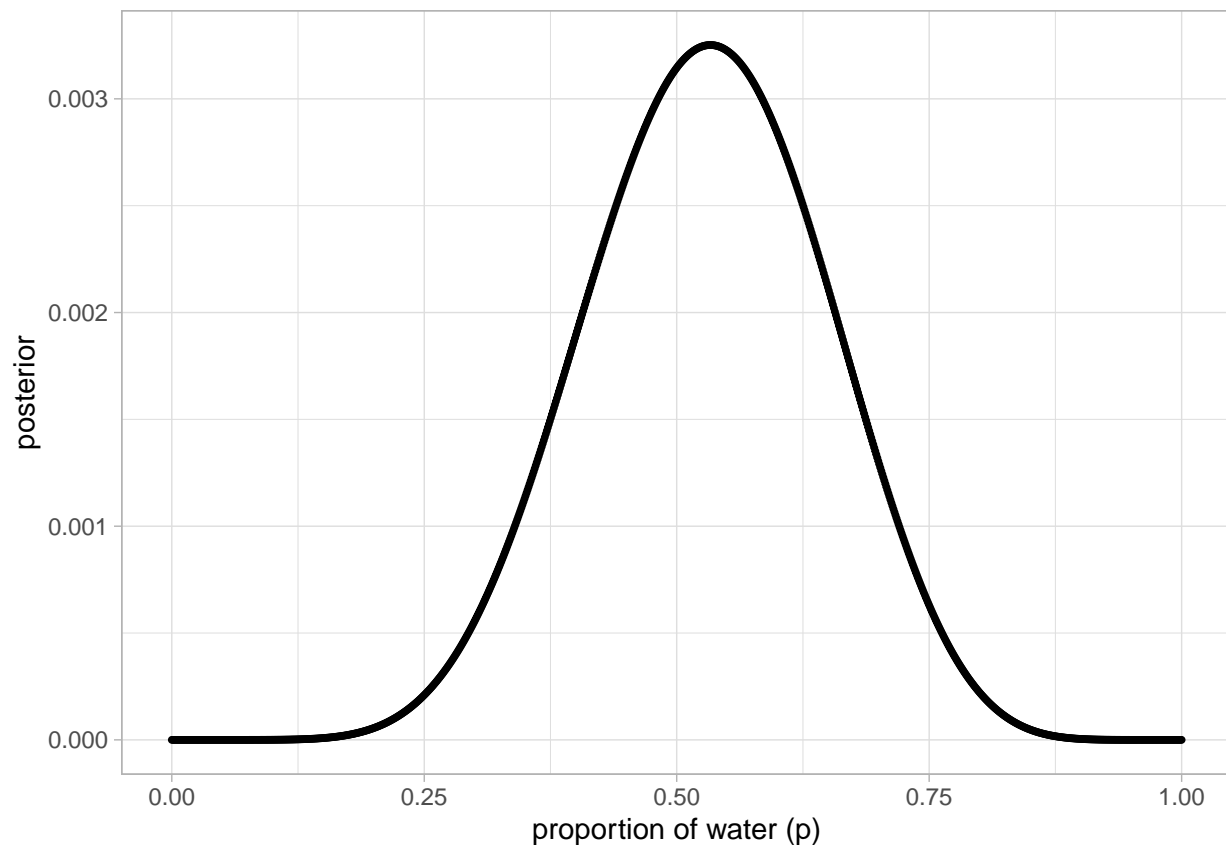
Medium

3M1. Suppose the globe tossing data had turned out to be 8 water in 15 tosses. Construct the posterior distribution, using grid approximation. Use the same flat prior as before.

```
n <- 1000
n_success <- 8
n_trials <- 15

data <- tibble(p_grid = seq(from = 0, to = 1, length.out = n),
               prior = 1) %>%
  mutate(likelihood = dbinom(n_success, size = n_trials, prob = p_grid),
         posterior = (likelihood * prior) / sum(likelihood * prior))

ggplot(data, aes(x = p_grid, y = posterior)) +
  geom_point(size = 0.7, show.legend = F) +
  geom_line(show.legend = F) +
  xlab("proportion of water (p)") +
  theme_light()
```

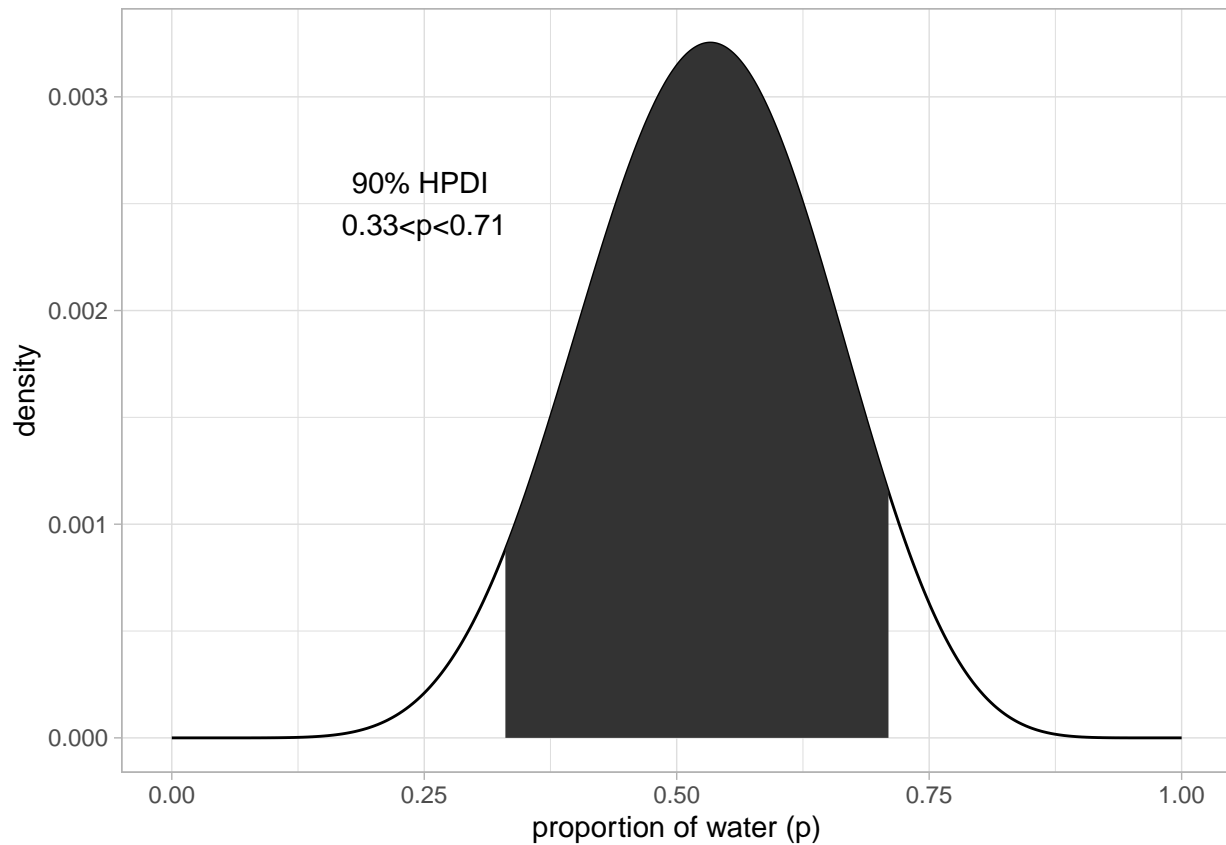


3M2. Draw 10,000 samples from the grid approximation from above. Then use the samples to calculate the 90% HPDI for p .

```
samples <- sample_n(data, size = 1e4, weight = posterior, replace = TRUE)

HPDI(samples$p_grid, prob = 0.9) %>%
  round(2) %>% c(lower, upper)

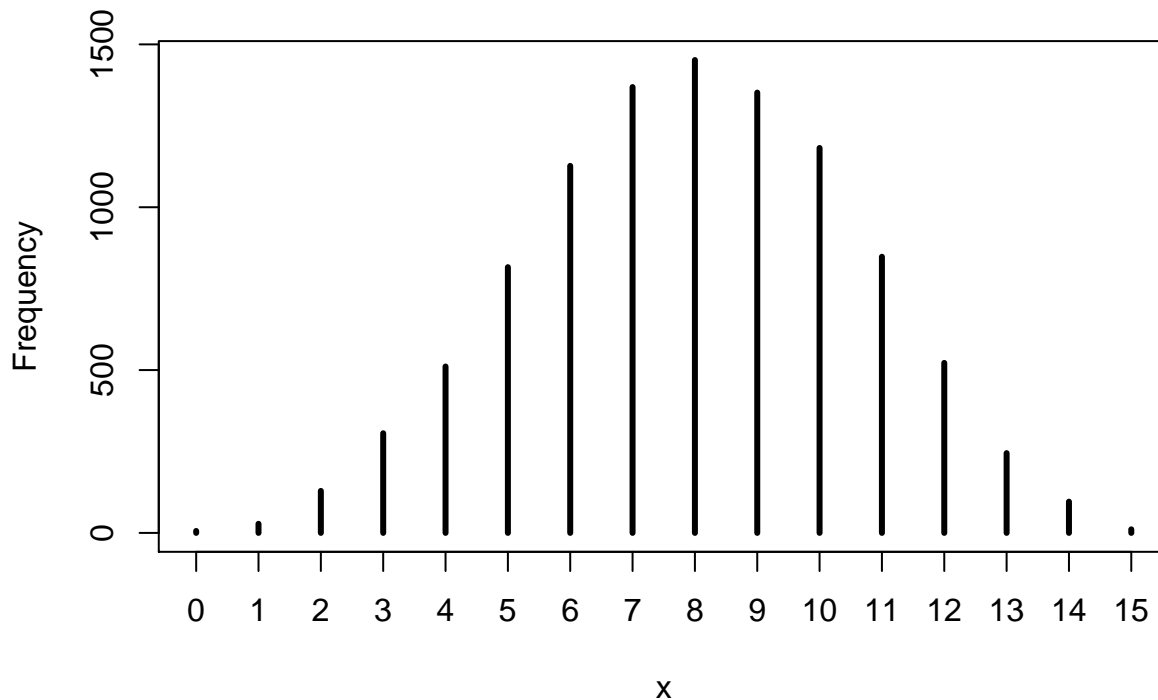
data %>%
  ggplot(aes(x = p_grid)) +
  geom_line(aes(y = posterior)) +
  geom_ribbon(data = data %>% filter(p_grid > lower, p_grid < upper),
            aes(ymin = 0, ymax = posterior)) +
  annotate(geom = "text",
          x = 0.25, y = 0.0025,
          label = paste0("90% HPDI \n",
                        lower, "<p<", upper)) +
  labs(x = "proportion of water (p)",
       y = "density") +
  theme_light()
```



3M3. Construct a posterior predictive check for this model and data. This means simulate the distribution of samples, averaging over the posterior uncertainty in p . What is the probability of observing 8 water in 15 tosses?

```
water <- rbinom(1e4, size = 15, prob = samples$p_grid)

simplehist(water)
```

```
cat("Probability of observing 8 water in 15 tosses is ",
    mean(water==8)*100, "%", sep = "")
```

Probability of observing 8 water in 15 tosses is 14.52%

3M4. Using the posterior distribution constructed from the new (8/15) data, now calculate the probability of observing 6 water in 9 tosses.

```
water <- rbinom(1e4, size = 9, prob = samples$p_grid)
```

```
cat("Probability of observing 6 water in 9 tosses is ",
    mean(water==6)*100, "%", sep = "")
```

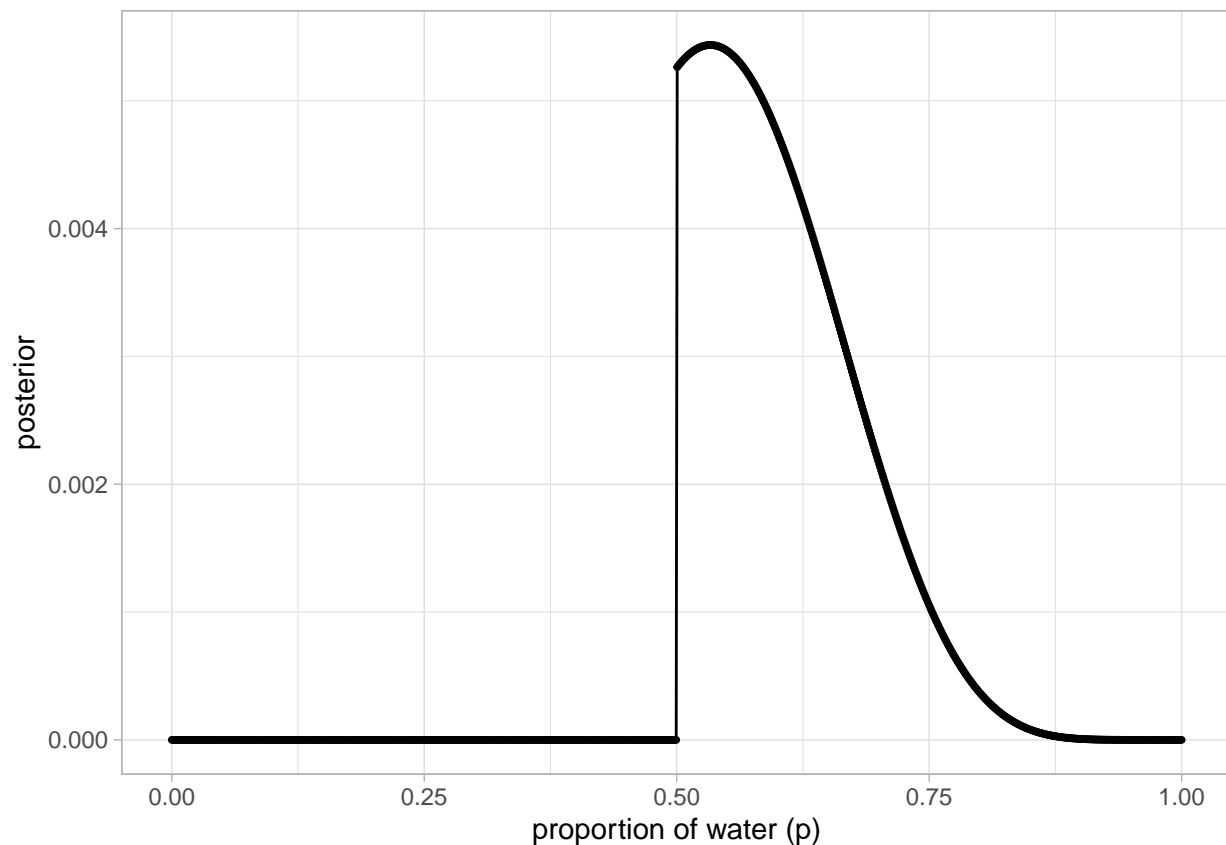
Probability of observing 6 water in 9 tosses is 18.16%

3M5. Start over at 3M1, but now use a prior that is zero below $p=0.5$ and a constant above $p=0.5$. This corresponds to prior information that a majority of the Earth's surface is water. Repeat each problem above and compare the inferences. What difference does the better prior make? If it helps, compare inferences (using both priors) to the true value $p = 0.7$.

```
n <- 1000
n_success <- 8
n_trials <- 15

data <- tibble(p_grid = seq(from = 0, to = 1, length.out = n),
               prior = if_else(p_grid < 0.5, 0, 1)) %>%
  mutate(likelihood = dbinom(n_success, size = n_trials, prob = p_grid),
         posterior = (likelihood * prior) / sum(likelihood * prior))

ggplot(data, aes(x = p_grid, y = posterior)) +
  geom_point(size = 0.7, show.legend = F) +
  geom_line(show.legend = F) +
  xlab("proportion of water (p)") +
  theme_light()
```



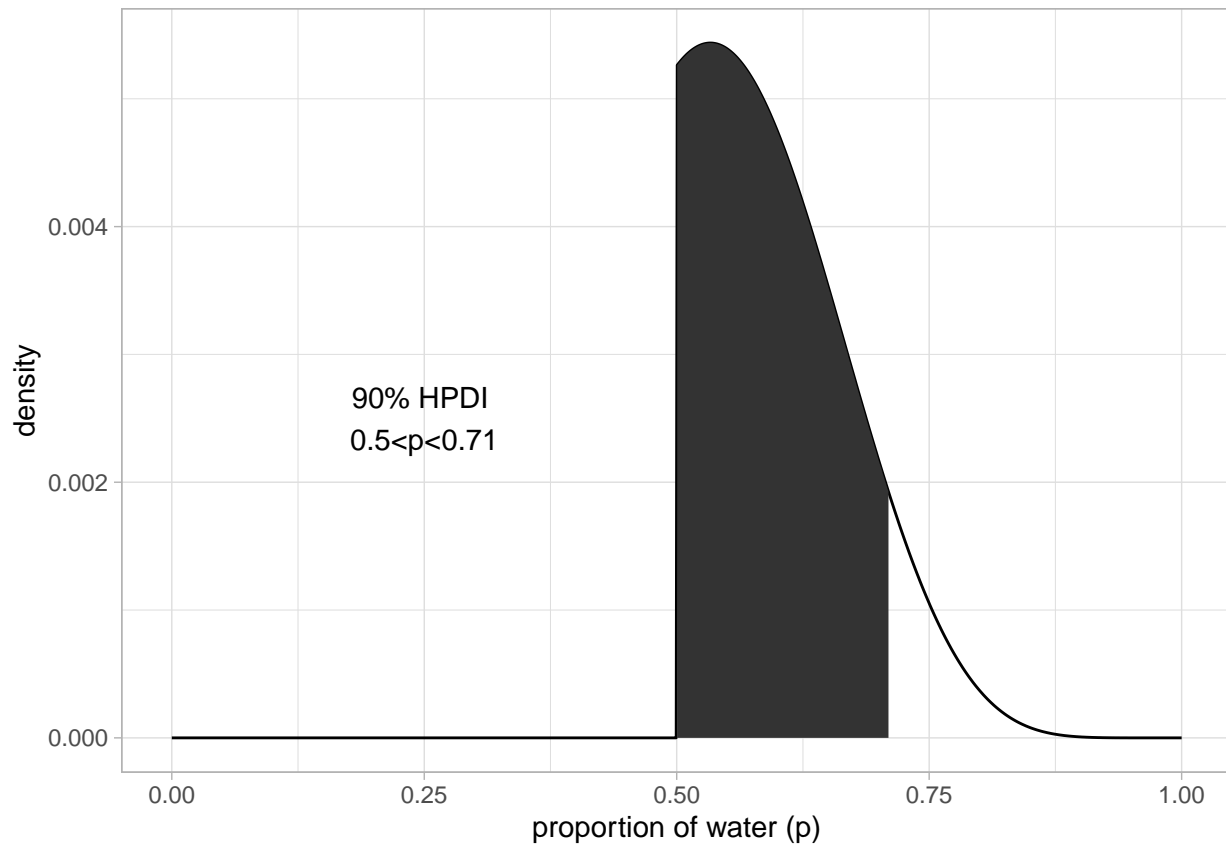
```

samples<- sample_n(data, size = 1e4, weight = posterior, replace = TRUE)

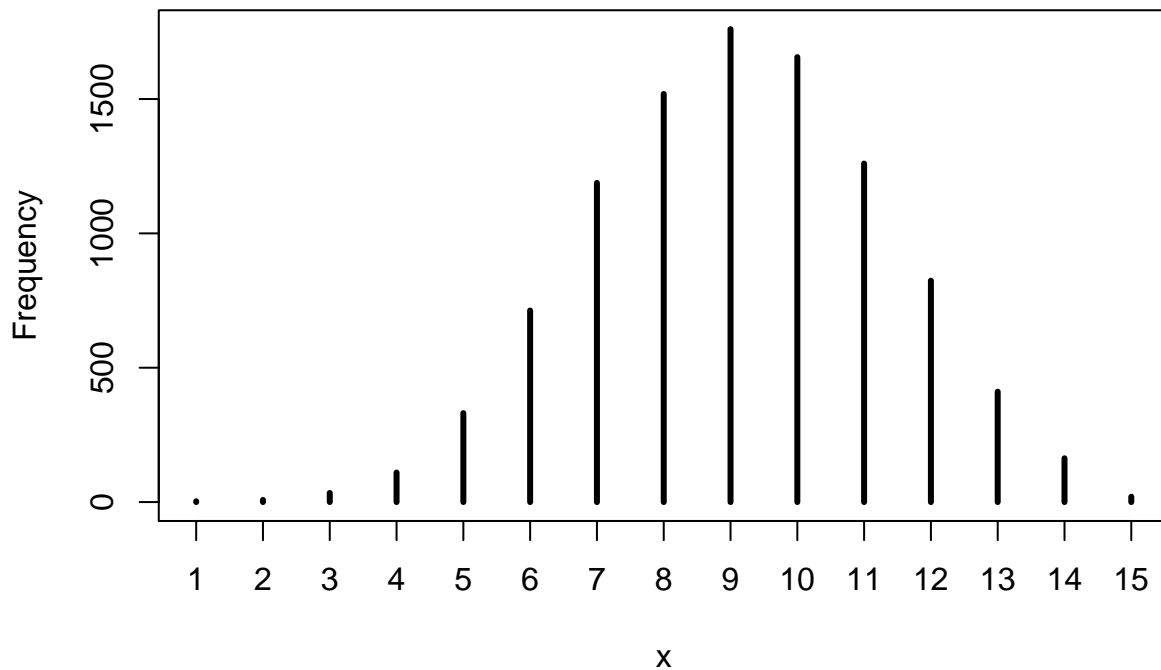
HPDI(samples$p_grid, prob = 0.9) %>%
  round(2) %>% c(lower, upper)

data %>%
  ggplot(aes(x = p_grid)) +
  geom_line(aes(y = posterior)) +
  geom_ribbon(data = data %>% filter(p_grid > lower, p_grid < upper),
            aes(ymin = 0, ymax = posterior)) +
  annotate(geom = "text",
          x = 0.25, y = 0.0025,
          label = paste0("90% HPDI \n",
                        lower, "<p<", upper)) +
  labs(x = "proportion of water (p)",
       y = "density") +
  theme_light()

```



```
water <- rbinom(1e4, size = 15, prob = samples$p_grid)
simplehist(water)
```



```
cat("Probability of observing 8 water in 15 tosses is ",
    mean(water==8)*100, "%", sep = "")
```

```
## Probability of observing 8 water in 15 tosses is 15.19%
```

```
water <- rbinom(1e4, size = 9, prob = samples$p_grid)
```

```
cat("Probability of observing 6 water in 9 tosses is ",  
    mean(water==6)*100, "%", sep = "")
```

```
## Probability of observing 6 water in 9 tosses is 23.12%
```

3M6. Suppose you want to estimate the Earth's proportion of water very precisely. Specifically, you want the 99% percentile interval of the posterior distribution of p to be only 0.05 wide. This means the distance between the upper and lower bound of the interval should be 0.05. How many times will you have to toss the globe to do this?

```
n_trials <- 10  
n <- 1000  
p <- 0.7  
  
repeat {  
  n_trials <- n_trials+1  
  n_success <- floor(p*n_trials)  
  
  data <- tibble(p_grid = seq(from = 0, to = 1, length.out = n),  
                 prior = p) %>%  
    mutate(likelihood = dbinom(n_success, size = n_trials, prob = p_grid),  
           posterior = (likelihood * prior) / sum(likelihood * prior))  
  
  samples<- sample_n(data, size = 1e4, weight = posterior, replace = TRUE)  
  
  PI(samples$p_grid, prob = 0.9) %>%  
    round(2) %>% c(lower, upper)  
  
  if(upper-lower <= 0.05) {  
    cat("after ", n_trials,  
        " trials, the difference between the bounds is ",  
        upper - lower, sep = "")  
    break  
  }  
}
```

```
## after 730 trials, the difference between the bounds is 0.05
```

Hard

Introduction. The practice problems here all use the data below. These data indicate the gender (male=1, female=0) of officially reported first and second born children in 100 two-child families.

```
birth1 <- c(1,0,0,0,1,1,0,1,0,1,0,0,1,1,0,1,1,0,0,0,1,0,0,0,1,0,  
0,0,0,1,1,1,0,1,0,1,1,1,0,1,0,1,1,0,1,0,0,1,1,0,1,0,0,0,0,0,0,  
1,1,0,1,0,0,1,0,0,0,1,0,0,1,1,1,1,0,1,0,1,1,1,1,1,0,0,1,0,1,1,0,  
1,0,1,1,1,0,1,1,1,1)
```

```
birth2 <- c(0,1,0,1,0,1,1,1,0,0,1,1,1,1,0,0,1,1,1,0,0,1,1,1,0,  
1,1,1,0,1,1,1,0,1,0,0,1,1,1,1,0,0,1,0,1,1,1,1,1,1,1,1,1,1,1,  
1,1,1,0,1,1,0,1,1,0,1,1,1,0,0,0,0,0,0,1,0,0,0,1,1,0,0,1,0,0,1,1,  
0,0,0,1,1,1,0,0,0,0)
```

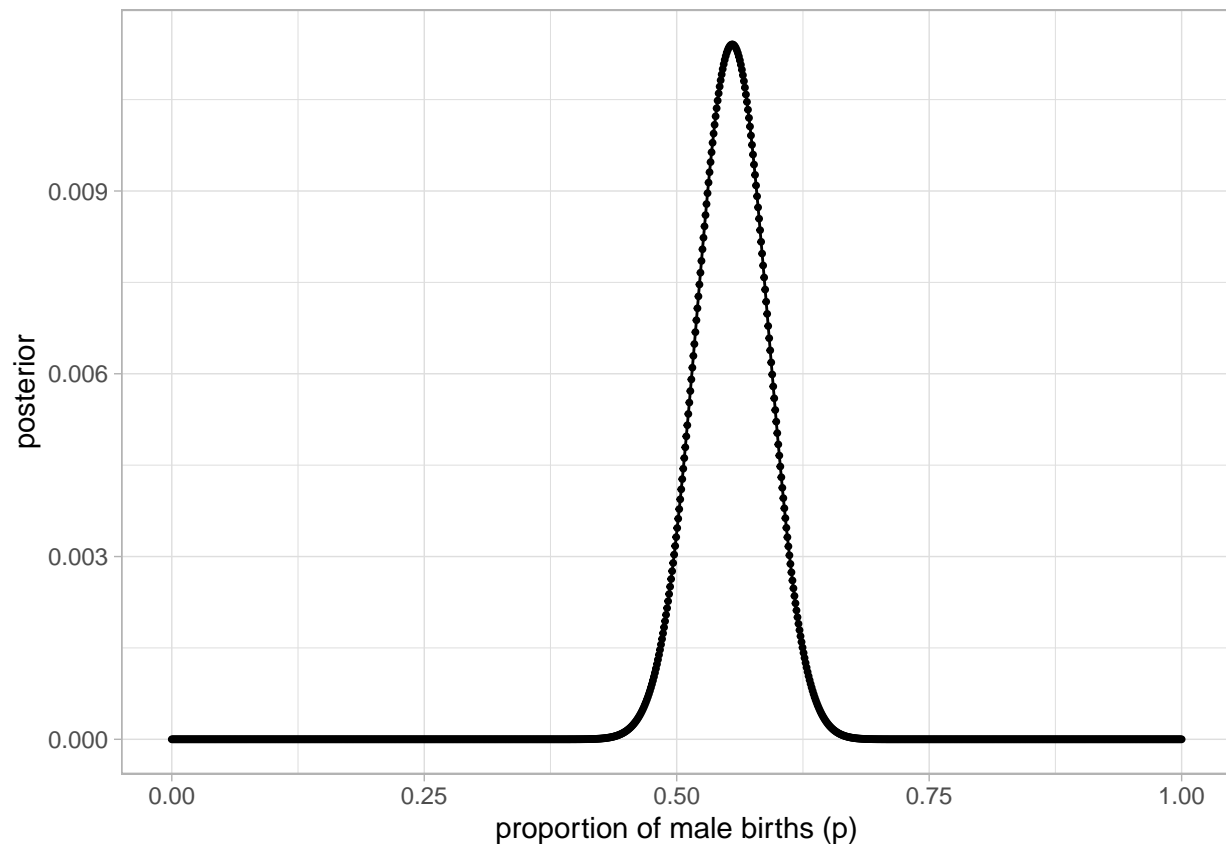
```
births <- c(birth1, birth2)
```

3H1. Using grid approximation, compute the posterior distribution for the probability of a birth being a boy. Assume a uniform prior probability. Which parameter value maximizes the posterior probability?

```
n <- 1000
n_success <- sum(births)
n_trials <- length(births)

data <- tibble(p_grid = seq(from = 0, to = 1, length.out = n),
               prior = 0.5) %>%
  mutate(likelihood = dbinom(n_success, size = n_trials, prob = p_grid),
         posterior = (likelihood * prior) / sum(likelihood * prior))

ggplot(data, aes(x = p_grid, y = posterior)) +
  geom_point(size = 0.7, show.legend = F) +
  geom_line(show.legend = F) +
  xlab("proportion of male births (p)") +
  theme_light()
```



```
data$p_grid[which.max(data$posterior)]
```

```
## [1] 0.5545546
```

3H2. Using the sample function, draw 10,000 random parameter values from the posterior distribution you calculated above. Use these samples to estimate the 50%, 89%, and 97% highest posterior density intervals.

```

samples <- sample_n(data, size = 1e4, weight = posterior, replace = TRUE)

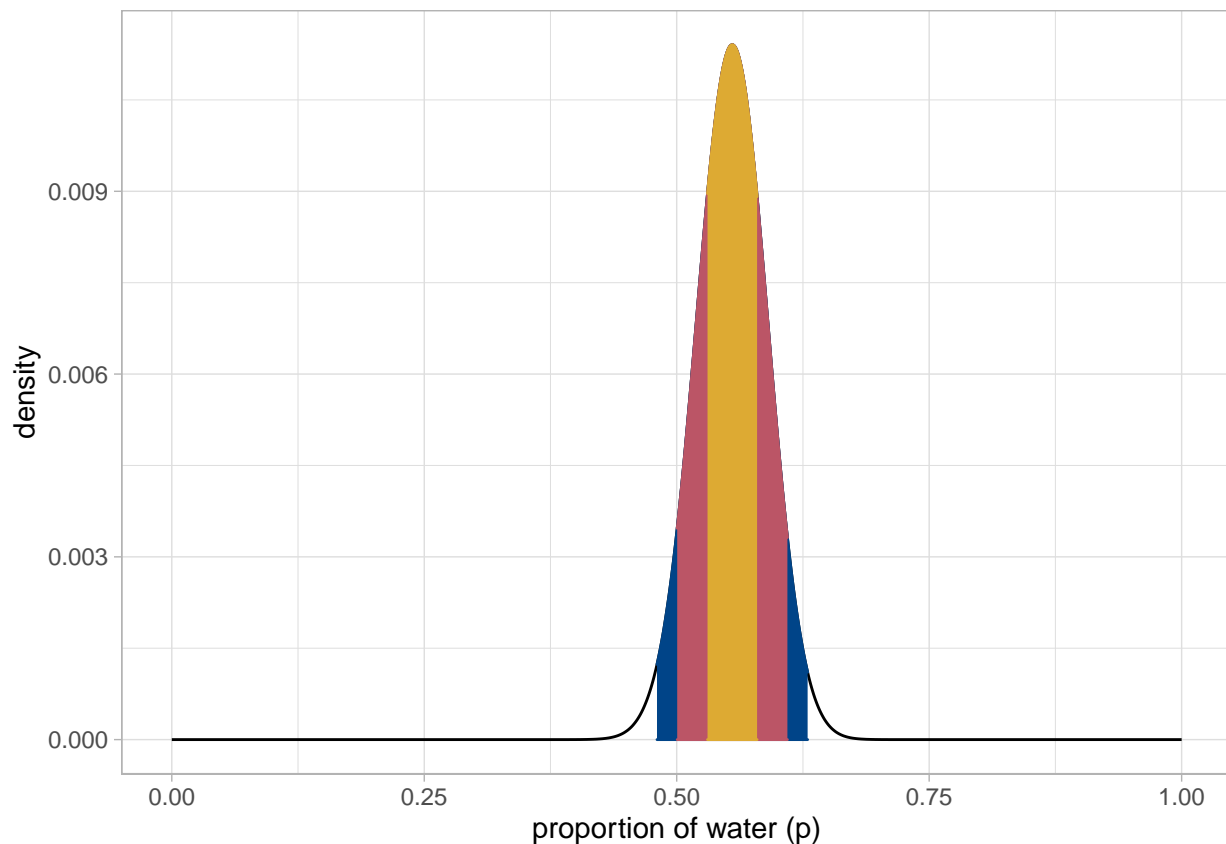
HPDI(samples$p_grid, prob = 0.5) %>%
  round(2) %>% c(lower_5, upper_5)

HPDI(samples$p_grid, prob = 0.89) %>%
  round(2) %>% c(lower_8, upper_8)

HPDI(samples$p_grid, prob = 0.97) %>%
  round(2) %>% c(lower_9, upper_9)

data %>%
  ggplot(aes(x = p_grid)) +
  geom_line(aes(y = posterior)) +
  geom_ribbon(data = data %>% filter(p_grid > lower_9, p_grid < upper_9),
            aes(ymin = 0, ymax = posterior), color = "#004488", fill = "#004488") +
  geom_ribbon(data = data %>% filter(p_grid > lower_8, p_grid < upper_8),
            aes(ymin = 0, ymax = posterior), color = "#BB5566", fill = "#BB5566") +
  geom_ribbon(data = data %>% filter(p_grid > lower_5, p_grid < upper_5),
            aes(ymin = 0, ymax = posterior), color = "#DDAA33", fill = "#DDAA33") +
  labs(x = "proportion of water (p)",
       y = "density") +
  theme_light()

```

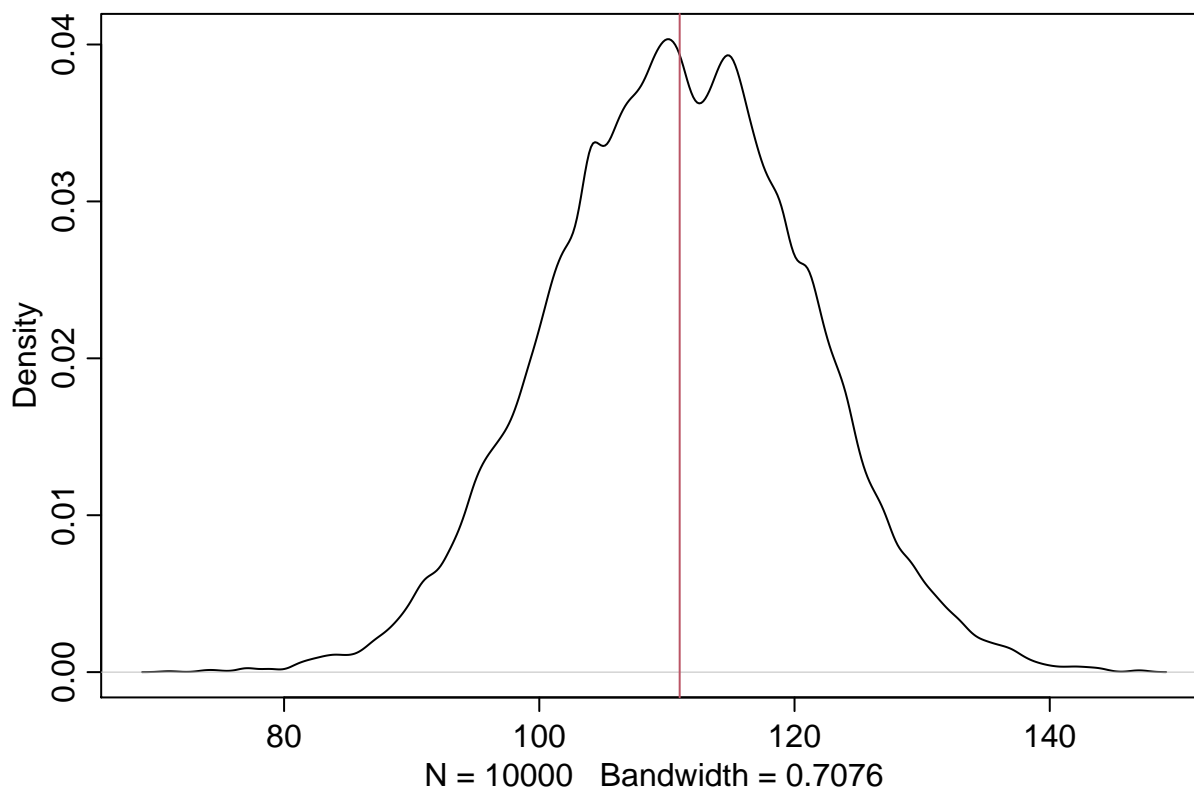


3H3. Use `rbinom` to simulate 10,000 replicates of 200 births. You should end up with 10,000 numbers, each one a count of boys out of 200 births. Compare the distribution of predicted numbers of boys to the actual count in the data (111 boys out of 200 births). There are many good ways to visualize the simulations, but

the `dens` command (part of the `rethinking` package) is probably the easiest way in this case. Does it look like the model fits the data well? That is, does the distribution of predictions include the actual observation as a central, likely outcome?

```
boys <- rbinom(1e4, size = 200, prob = samples$p_grid)
```

```
dens(boys)
abline(v = sum(births), col = "#BB5566" )
```



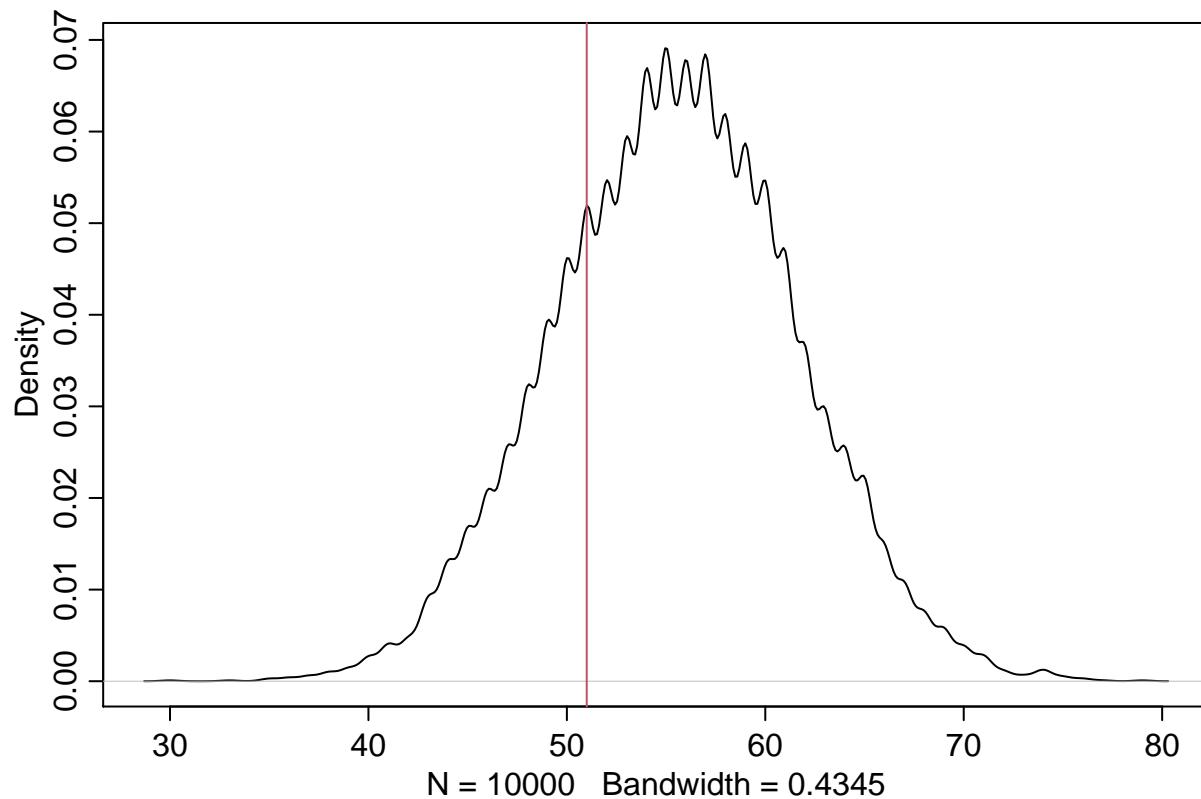
```
mean(boys) %>% floor()
```

```
## [1] 110
```

3H4. Now compare 10,000 counts of boys from 100 simulated first-borns only to the number of boys in the first births, *birth1*. How does the model look in this light?

```
boys <- rbinom(1e4, size = 100, prob = samples$p_grid)
```

```
dens(boys)
abline(v = sum(birth1), col = "#BB5566" )
```



```
cat("The observed value of ", sum(birth1),
    " in 100 first borns (births1) is not represented well with the model
    since simulated results have a median of ", median(boys),
    " and a mean of ", mean(boys) %>% floor(), " .", sep = "")
```

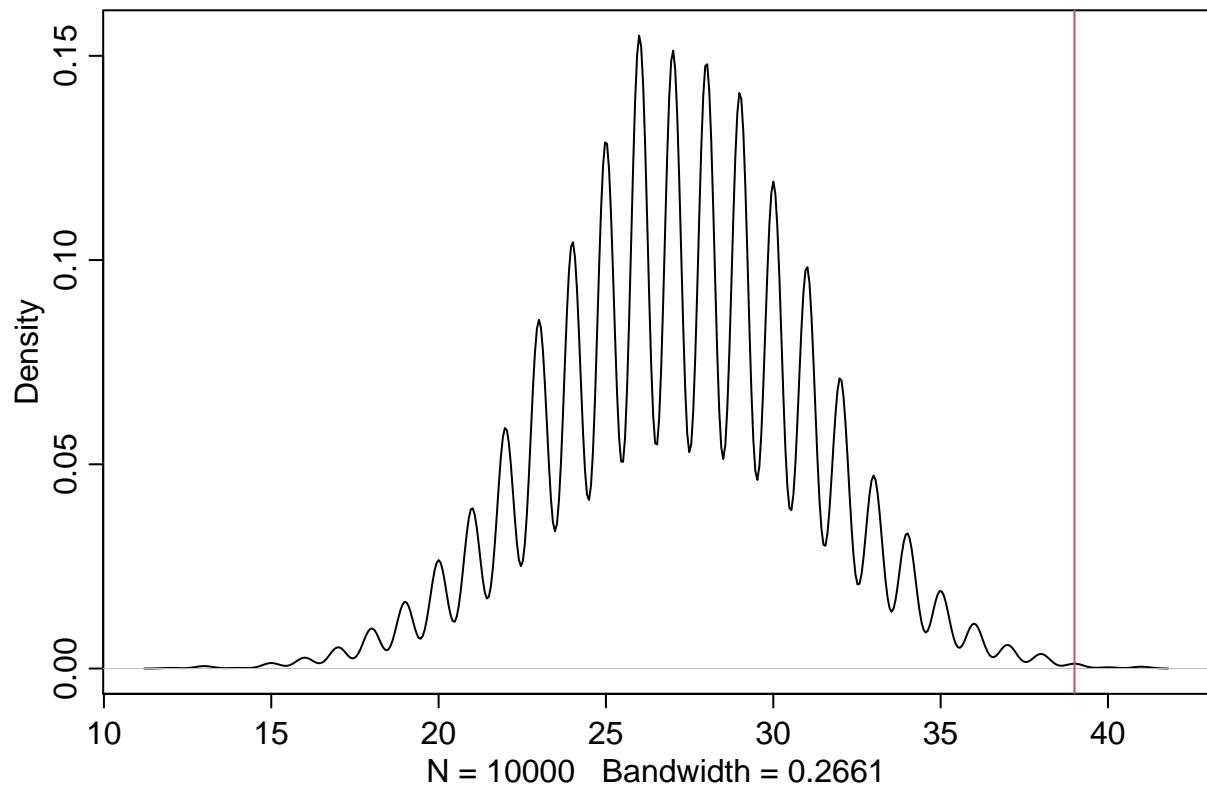
```
## The observed value of 51 in 100 first borns (births1) is not represented well with the model
##     since simulated results have a median of 56 and a mean of 55 .
```

3H5. The model assumes that sex of first and second births are independent. To check this assumption, focus now on second births that followed female first-borns. Compare 10,000 simulated counts of boys to only those second births that followed girls. To do this correctly, you need to count the number of first-borns who were girls and simulate that many births, 10,000 times. Compare the counts of boys in your simulations to the actual observed count of boys following girls. How does the model look in this light? Any guesses what is going on in these data?

```
boys_after_girls <- birth2[birth1 == 0]

boys <- rbinom(1e4, size = length(boys_after_girls), prob = samples$p_grid)

dens(boys)
abline(v = sum(boys_after_girls), col = "#BB5566")
```

More at github.com/antoniojurlina/statistical_rethinking

sessionInfo()

```
## R version 3.6.2 (2019-12-12)
## Platform: x86_64-apple-darwin15.6.0 (64-bit)
## Running under: macOS Sierra 10.12.6
##
## Matrix products: default
## BLAS: /Library/Frameworks/R.framework/Versions/3.6/Resources/lib/libRblas.0.dylib
## LAPACK: /Library/Frameworks/R.framework/Versions/3.6/Resources/lib/libRlapack.dylib
##
## locale:
## [1] en_US.UTF-8/en_US.UTF-8/en_US.UTF-8/C/en_US.UTF-8/en_US.UTF-8
##
## attached base packages:
## [1] parallel stats graphics grDevices utils datasets methods
## [8] base
##
## other attached packages:
## [1] zeallot_0.1.0 patchwork_1.0.1 rethinking_2.11
## [4] rstan_2.21.1 StanHeaders_2.21.0-5 knitr_1.28
## [7] here_0.1 ggthemes_4.2.0 forcats_0.5.0
## [10] stringr_1.4.0 dplyr_0.8.5 purrr_0.3.4
## [13] readr_1.3.1 tidyr_1.0.3 tibble_3.0.3
## [16] ggplot2_3.3.2 tidyverse_1.3.0
##
## loaded via a namespace (and not attached):
## [1] httr_1.4.1 jsonlite_1.7.0 modelr_0.1.7 RcppParallel_5.0.2
## [5] assertthat_0.2.1 stats4_3.6.2 cellranger_1.1.0 yaml_2.2.1
## [9] pillar_1.4.6 backports_1.1.8 lattice_0.20-41 glue_1.4.1
## [13] digest_0.6.25 rvest_0.3.5 colorspace_1.4-1 htmltools_0.4.0
## [17] pkgconfig_2.0.3 broom_0.5.6 haven_2.2.0 mvtnorm_1.1-1
## [21] scales_1.1.1 processx_3.4.3 farver_2.0.3 generics_0.0.2
## [25] ellipsis_0.3.1 withr_2.2.0 cli_2.0.2 magrittr_1.5
## [29] crayon_1.3.4 readxl_1.3.1 evaluate_0.14 ps_1.3.3
## [33] fs_1.4.1 fansi_0.4.1 nlme_3.1-147 MASS_7.3-51.6
## [37] xml2_1.3.2 pkgbuild_1.1.0 tools_3.6.2 loo_2.3.1
## [41] prettyunits_1.1.1 hms_0.5.3 lifecycle_0.2.0 matrixStats_0.56.0
## [45] V8_3.2.0 munsell_0.5.0 reprex_0.3.0 callr_3.4.3
## [49] compiler_3.6.2 rlang_0.4.7 grid_3.6.2 rstudioapi_0.11
## [53] labeling_0.3 rmarkdown_2.1 gtable_0.3.0 codetools_0.2-16
## [57] inline_0.3.15 DBI_1.1.0 curl_4.3 R6_2.4.1
## [61] gridExtra_2.3 lubridate_1.7.8 rprojroot_1.3-2 shape_1.4.4
## [65] stringi_1.4.6 Rcpp_1.0.5 vctrs_0.3.2 dbplyr_1.4.3
## [69] tidyselect_1.1.0 xfun_0.13 coda_0.19-3
```