

Applied Bayesian modeling - HW2

Álvaro J. Castro Rivadeneira

September 16, 2022

Score: The maximum number of points in this HW is 12 points, with 3 points extra credit. For calculating a final HW grade, the points will be rescaled to a maximum score of $(12+3)/12 \cdot 100\% = 125\%$.

What to hand in: For any exercises using R, we need an Rmd and a knitted pdf. You can add the answers to other exercises in the same markdown file or create them in a different way (as long as it's legible). If you create them outside markdown, please combine all answers into one pdf.

Part 1 - based on module 4

Exercise 1: Derivation of a posterior using Bayes' rule [3 pts]

This exercise is about the material in module 4.

Obtain $p(\mu|\mathbf{y})$ when prior and data are defined as follows:

$$\begin{aligned} y_i|\theta_i, \sigma^2 &\sim N(\theta_i, \sigma^2) \text{ (independent), for } i = 1, 2, \dots, n; \\ \theta_i &= \mu + r_i; \\ \mu &\sim N(m_0, s_{\mu 0}^2); \end{aligned}$$

where r_i refers to a known/fixed reporting error (that varies across observations) and σ^2 is known.

Note that you can treat the r_i as fixed, not random.

(A more realistic but somewhat more difficult exercise would be to have r_i be random, and assume that they are normally distributed and independent of the y_i 's).

Answer:

If every y_i is normally distributed around θ_i , then \bar{y} is also normally distributed around a mean $\bar{\theta}$, with a variance $\frac{\sigma^2}{n}$.

As was explained in class, we can define a variable $z_i = y_i - r_i$ such that:

$$\bar{z} = \frac{1}{n} \cdot \sum_{i=1}^n (y_i - r_i) = \frac{1}{n} \cdot \sum_{i=1}^n (y_i) - \frac{1}{n} \cdot \sum_{i=1}^n (r_i) = \bar{y} - \bar{r}$$

Given this, and that $\mu = \theta_i - r_i$, and that the variance only depends on θ_i since r_i is fixed, then:

$$z_i|\mu, \sigma^2 \sim N(\mu, \sigma^2)$$

So, going back to Bayes' rule, we can derive:

$$\begin{aligned} p(\mu|z) &\propto p(\mu)p(z|\mu) \\ p(\mu|z) &\propto \exp\left(\frac{-1}{2s_{\mu 0}^2}(\mu - m_0)^2\right) \cdot \exp\left(\frac{-n}{2\sigma^2}(\bar{z} - \mu)^2\right) \end{aligned}$$

$$p(\mu|z) \propto \exp\left(-\frac{1}{2}f(\mu)\right)$$

where

$$\begin{aligned} f(\mu) &= 1/s_{\mu_0}^2 (\mu^2 - 2 \cdot m_0 \cdot \mu) + n/\sigma^2 (\mu^2 - 2 \cdot \bar{z} \cdot \mu) \\ f(\mu) &\propto (1/s_{\mu_0}^2 + n/\sigma^2) \mu^2 - 2 \cdot (1/s_{\mu_0}^2 m_0 + n/\sigma^2 \bar{z}) \cdot \mu \end{aligned}$$

So that:

$$\mu|z \sim N\left(\frac{m_0/s_{\mu_0}^2 + n \cdot \bar{z}/\sigma^2}{1/s_{\mu_0}^2 + n/\sigma^2}, \frac{1}{1/s_{\mu_0}^2 + n/\sigma^2}\right)$$

Substituting $\bar{z} = \bar{y} - \bar{r}$ we get:

$$p(\mu|(y-r)) \propto p(\mu)p((y-r)|\mu)$$

Since we assume that r is fixed, then $p(\mu|z) = p(\mu|y, r) = p(\mu|y)$

$$\begin{aligned} p(\mu|z) &= p(\mu|y) \propto \exp\left(\frac{-1}{2s_{\mu_0}^2}(\mu - m_0)^2\right) \cdot \exp\left(\frac{-n}{2\sigma^2}((\bar{y} - \bar{r}) - \mu)^2\right) \\ p(\mu|y) &\propto \exp\left(-\frac{1}{2}f(\mu)\right) \end{aligned}$$

where

$$\begin{aligned} f(\mu) &= 1/s_{\mu_0}^2 (\mu^2 - 2 \cdot m_0 \cdot \mu) + n/\sigma^2 (\mu^2 - 2 \cdot (\bar{y} - \bar{r}) \cdot \mu) \\ f(\mu) &\propto (1/s_{\mu_0}^2 + n/\sigma^2) \mu^2 - 2 \cdot (1/s_{\mu_0}^2 m_0 + n/\sigma^2 (\bar{y} - \bar{r})) \cdot \mu \end{aligned}$$

So that:

$$\mu|y \sim N\left(\frac{m_0/s_{\mu_0}^2 + n \cdot (\bar{y} - \bar{r})/\sigma^2}{1/s_{\mu_0}^2 + n/\sigma^2}, \frac{1}{1/s_{\mu_0}^2 + n/\sigma^2}\right)$$

Exercise 2: Interpretation of the role of prior information for estimating IQ scores

This exercise is about the material in module 4, based on Hoff 5.4 (but using a different prior on μ).

Background: Scoring on IQ tests is designed to produce a normal distribution with a mean of 100 and a standard deviation of 15 (a variance of 225) when applied to the general population.

Suppose we are to sample n individuals from a particular town in the United States and then estimate μ , the town-specific mean IQ score, based on the sample of size n .

Let y_i denote the IQ score for the i -th person in the town of interest, and assume $y_1, y_2, \dots, y_n | \mu, \sigma^2 \sim N(\mu, \sigma^2)$ (independent), with $\sigma = 15$. Suppose that $\bar{y} = 113$ and $n = 10$.

For Bayesian inference about μ , the following prior will be used:

$$\mu \sim N(m_0, s_{\mu_0}^2),$$

with $m_0 = 100$ and $s_{\mu_0} = \sigma = 15$ (based on the information about IQ scores).

Exercise 2a [5pts]

- (i) Given the information above, obtain the expression for the Bayesian point estimate of μ , $\hat{\mu}_{Bayes} = E(\mu|\mathbf{y})$, in terms of m_0 , n , \bar{y} , and/or σ (don't plug in values yet; note that the expression simplifies relative to the expression we obtained in module 4, using only a subset of these).

For the point estimate of μ , $\hat{\mu}_{Bayes} = E(\mu|\mathbf{y})$, we can use the mean of the posterior, which is:

$$\hat{\mu}_{Bayes} = \frac{m_0/s_{\mu_0}^2 + n \cdot \bar{y}/\sigma^2}{1/s_{\mu_0}^2 + n/\sigma^2}$$

Since $s_{\mu_0} = \sigma$, I can substitute the former with the latter to get:

$$\begin{aligned}\hat{\mu}_{Bayes} &= \frac{m_0/\sigma^2 + n \cdot \bar{y}/\sigma^2}{1/\sigma^2 + n/\sigma^2} = \frac{m_0 + n \cdot \bar{y}}{\sigma^2} \cdot \frac{\sigma^2}{(1+n)} \\ \hat{\mu}_{Bayes} &= \frac{m_0 + n \cdot \bar{y}}{(1+n)}\end{aligned}$$

Which is surprising insofar the variance is removed from the equation.

- (ii) Interpret the Bayesian point estimate as the weighted combination of prior information and data.

Since the variance is the same in the prior and the likelihood, then it doesn't contribute to weighting both of these, as it will give equal weight to both. Thus, the Bayesian mean becomes simply an average of all the data points including the prior mean as another data point. This can be seen more clearly if we assumed that $m_0 = \bar{y}$, so that the numerator above becomes $\bar{y} + n \cdot \bar{y} = (1+n) \cdot \bar{y}$. Thus, you would simply end up with \bar{y} as the point estimate. The weight of the prior depends on how many data points there are, and its weight is ultimately just that of another data point.

- (iii) Then calculate the value for $\hat{\mu}_{Bayes}$ given the information provided and construct a 95% credible interval for μ .

Plugging in the values, we get

$$\hat{\mu}_{Bayes} = \frac{100 + 10 \cdot 113}{(1 + 10)} = \frac{1230}{11} = 111.818$$

```
(mu.hat = (100 + (10*113)) / (1 + 10))
```

```
## [1] 111.8182
```

To get the 95% credible interval, I can use:

```
qnorm(c(0.025, 0.975), mean = mu.hat, sd = 15) # 95% quantile-based CI
```

```
## [1] 82.41872 141.21764
```

This means that my 95% credible interval ranges from 82.42 - 141.22, which notably includes this town's mean.

2b Extra credit [3pts]

We will now compare the sampling properties of the Bayes estimator for μ , $\hat{\mu}_{Bayes} = E(\mu|\mathbf{y})$, to the maximum likelihood estimator $\hat{\mu}_{MLE} = \bar{y}$. Sampling properties of estimators refer to their behavior under hypothetically repeatable surveys or experiments, summarized into bias and mean squared error, as explained below with μ^* referring to the (unknown) true value of μ :

- The bias of an estimator is defined as the difference between its expected value (w.r.t. data \mathbf{y}) and the (unknown) true value:

$$Bias(\hat{\mu}) = E_{\mathbf{y}}[\hat{\mu}|\mu^*] - \mu^*$$

(adding the subscript with the expectation here to make it clear that the expectation is wrt \mathbf{y} , to clarify that the data are what's random here; also adding conditioning on true value μ^*)

- Bias refers to how close the center of mass of the sampling distribution of an estimator is to the true value. An unbiased estimator is an estimator with zero bias, which sounds desirable. However, bias does not tell us how far away an estimate might be from the true value. For example, y_1 is an unbiased estimator of the population mean μ , but will generally be farther away from μ than \bar{y} .

- To evaluate how close an estimator $\hat{\mu}$ is likely to be to the true value μ^* , we might use the mean squared error (MSE). Letting $m = E_{\mathbf{y}}[\hat{\mu}|\mu^*]$, the MSE of an estimator $\hat{\mu}$ is

$$\begin{aligned} MSE[\hat{\mu}|\mu^*] &= E_{\mathbf{y}}[(\hat{\mu} - \mu^*)^2|\mu^*], \\ &= E_{\mathbf{y}}[(\hat{\mu} - m + m - \mu^*)^2|\mu^*], \\ &= E_{\mathbf{y}}[(\hat{\mu} - m)^2|\mu^*] + 2E_{\mathbf{y}}[(\hat{\mu} - m)(m - \mu^*)] + E[(m - \mu^*)^2|\mu^*], \\ &= E_{\mathbf{y}}[(\hat{\mu} - m)^2|\mu^*] + (m - \mu^*)^2, \\ &= Var[\hat{\mu}|\mu^*] + Bias(\hat{\mu})^2. \end{aligned}$$

This means that, before the data are gathered, the expected distance from the estimator to the true value depends on how close μ^* is to the center of the distribution of $\hat{\mu}$ (the bias), as well as how spread out the distribution is (the variance of $\hat{\mu}$).

Exercise Suppose that (unknown to us) the true mean IQ score μ^* in the town was quite high, $\mu^* = 112$. Calculate the bias, variance and MSE of the Bayes and ML estimators. Which estimator has a larger bias? Which estimator has a larger MSE?

Hint (based on a commonly made mistake): $E_{\mathbf{y}}(\hat{\mu}_{Bayes}|\mu^*) \neq \hat{\mu}_{Bayes}$, you need to get the expected value w.r.t. the data \mathbf{y} .

Optional: Consider

(i) plotting the sampling distributions for both the Bayes estimator as well as the MLE.

(ii) obtaining the Bayes and ML MSEs for sample sizes $n = 10$ to 1,000 and plotting the ratio (Bayes MSE)/(ML MSE) against sample size, to then interpret your findings.

Part 2 - based on module 5

Exercise 3: get stan going on your laptop [4pts]

The goal of this exercise is just to make sure you have stan (specifically, through rstan and brms) working on your laptop. To do so, please work through the Rmd file with module5, module5_sampling.Rmd.

For this HW Rmd, please add an example model fit using brms. This can be a model fit to radon data, copied from the module5 example, or a model fit from the brms examples. Also let us know if you successfully installed rstan and got it to work. If not, please explain the issue.

Answer:

I successfully installed `rstan` and got it to work. I worked through module 5, and to demonstrate it, here is a normal model fit from the `brms` examples with some slight modifications:

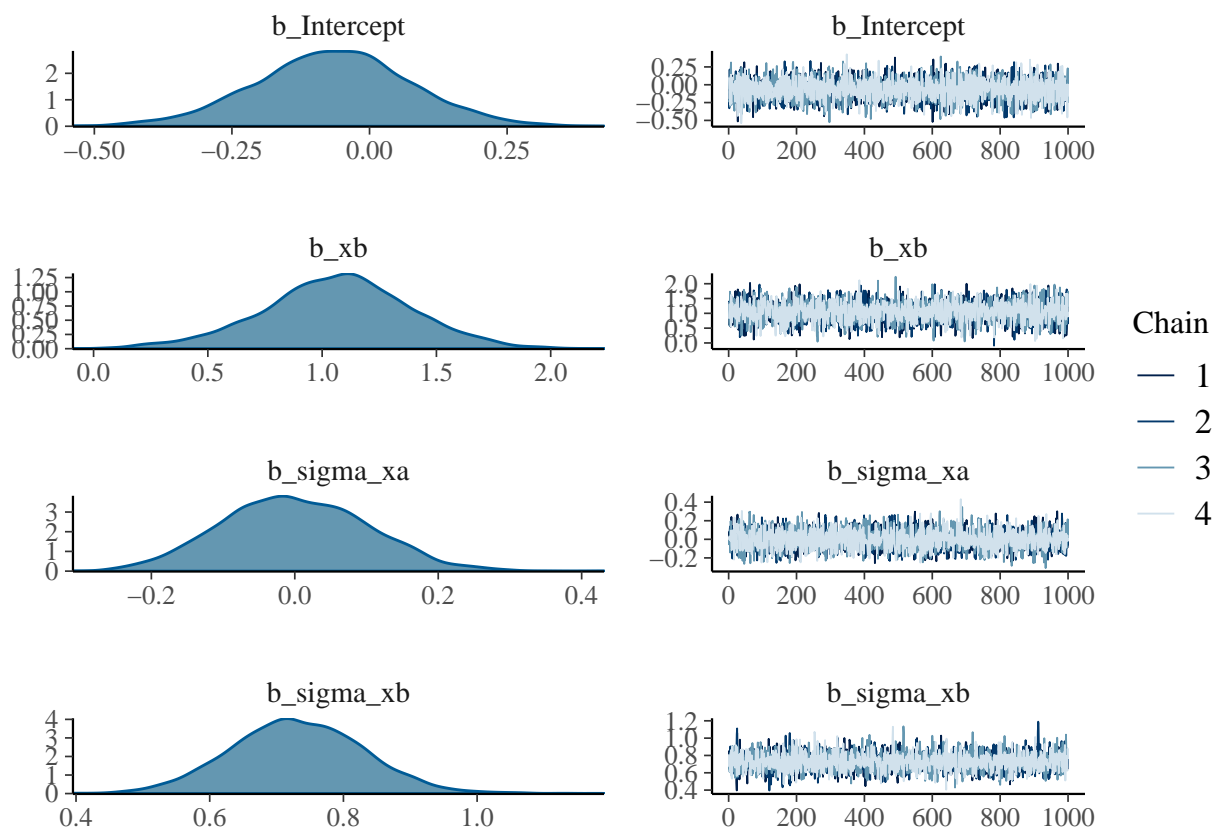
```
set.seed(1234)
mu_priorhw2 <- set_prior("normal(0,1)", class = "Intercept")
# Normal model with heterogeneous variances
data_het <- data.frame(
  y = c(rnorm(100), rnorm(100, 1, 2)),
  x = factor(rep(c("a", "b"), each = 100))
)
fit6 <- brm(bf(y ~ x, sigma ~ 0 + x), family = gaussian(), data = data_het,
  cores = getOption("mc.cores", 2),
  prior = c(mu_priorhw2),
  file = "output/exercise3"
)
summary(fit6)
```



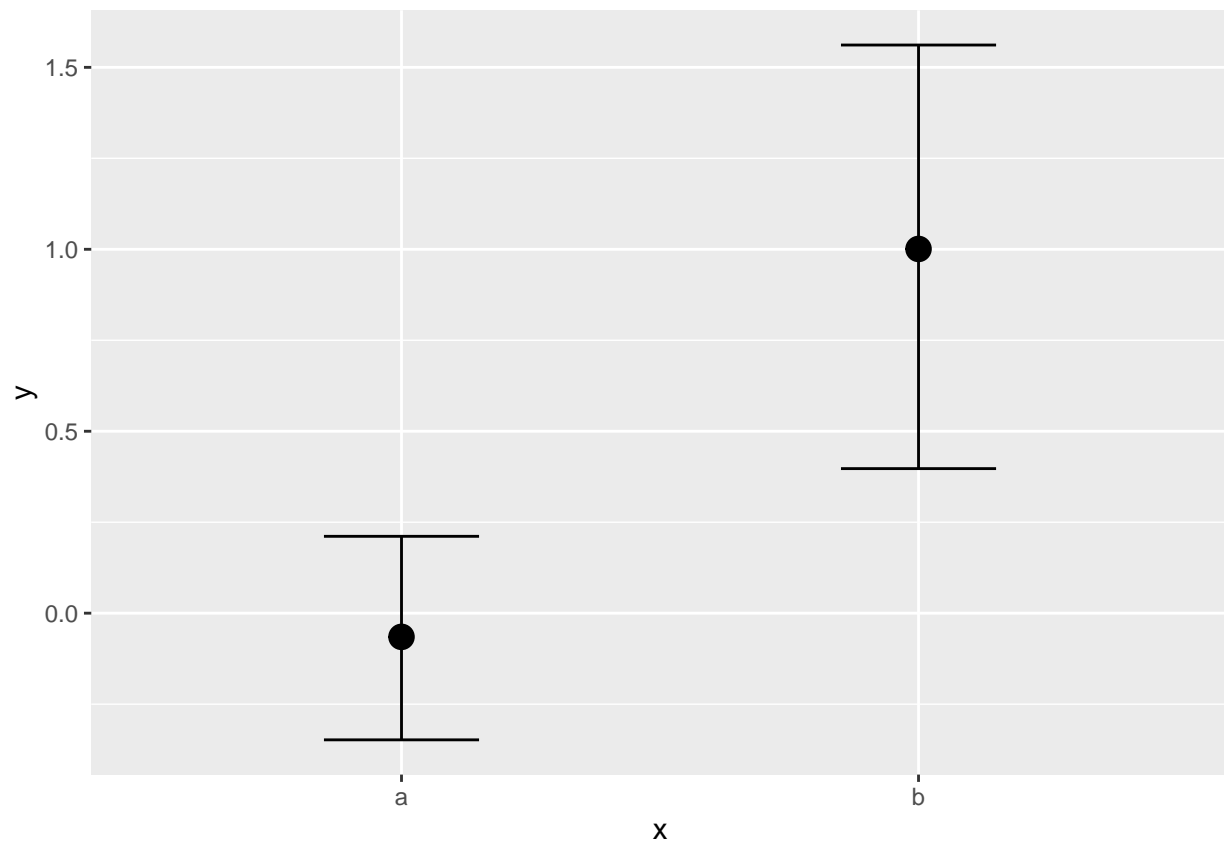
```
## Family: gaussian
## Links: mu = identity; sigma = log
## Formula: y ~ x
##          sigma ~ 0 + x
## Data: data_het (Number of observations: 100)
## Draws: 4 chains, each with iter = 2000; warmup = 1000; thin = 1;
##          total post-warmup draws = 4000
##
## Population-Level Effects:
##           Estimate Est.Error 1-95% CI u-95% CI Rhat Bulk_ESS Tail_ESS
## Intercept    -0.07     0.14   -0.35    0.21 1.00     6040     3168
## xb             1.06     0.32    0.39    1.69 1.00     3061     2784
## sigma_xa       0.00     0.10   -0.19    0.20 1.00     3380     2542
## sigma_xb       0.73     0.10    0.54    0.93 1.00     3428     2960
##
## Draws were sampled using sampling(NUTS). For each parameter, Bulk_ESS
## and Tail_ESS are effective sample size measures, and Rhat is the potential
## scale reduction factor on split chains (at convergence, Rhat = 1).
```



```
plot(fit6)
```



```
conditional_effects(fit6)
```



```
# extract estimated residual SDs of both groups
sigmas <- exp(as.data.frame(fit6, variable = "^b_sigma_", regex = TRUE))
ggplot(stack(sigmas), aes(values)) +
  geom_density(aes(fill = ind))
```

