

# Homework 3: Bayesian Analysis

Harvard CS 109B, Spring 2018

Feb 20, 2018

Homework 3 is due March 5, 2018 11:59 PM

## LDA & Bayes

In the first part of this assignment, you will be working with text from @realDonaldTrump Twitter. The text was taken from all tweets Donald Trump sent between 01/19/2016 and 01/19/2018. The goal is to use Latent Dirichlet Allocation in order to model the topics that the president tweeted about during this time.

In the second part of this assignment, you are provided with data sets *dataset-2-train.txt* and *dataset-2-test.txt* containing details of contraceptive usage by 1934 Bangladeshi women. There are four attributes for each woman, along with a label indicating if she uses contraceptives. The attributes include:

- district: identifying code for the district the woman lives in
- urban: type of region of residence
- living.children: number of living children
- age-mean: age of women (in years, centred around mean)

The women are grouped into 60 districts. The task is to build a classification model that can predict if a given woman uses contraceptives.

### 1. Data Preparation

The tweet data is provided for you as *trump-tibble.csv*. After you read the data into R, you'll see that there are only two columns: *document* and *text*. The *document* column contains the date and time of the tweet, and the *text* column contains the actual text of the tweet. Before you begin, you'll want to cast the columns as characters rather than factors. You can do this with the following code:

```
# read in trump tibble data
trump_tibble = read.csv("data/trump_tibble.csv")
# cast factors to characters
trump_tibble$document <- as.character(trump_tibble$document)
trump_tibble$text <- as.character(trump_tibble$text)
```

The following libraries will be of use for this problem:

```
# load libraries
library(topicmodels) #topic modeling functions
library(stringr) #common string functions
library(tidytext) #tidy text analysis
suppressMessages(library(tidyverse)) #data manipulation and visualization
#messages give R markdown compile error so we need to suppress it
```

```
## Source topicmodels2LDavis & optimal_k functions
```

```
invisible(lapply(file.path("https://raw.githubusercontent.com/trinker/topicmodels_learning/master/functions",
c("topicmodels2LDavis.R", "optimal_k.R")),
devtools::source_url))
```

```
## SHA-1 hash of file is 5ac52af21ce36dfe8f529b4fe77568ced9307cf0
```

```
## SHA-1 hash of file is 7f0ab64a94948c8b60ba29dddf799e3f6c423435
```

```
## Loading required package: pacman
```

(a) Use the `unnest-tokens` function to extract words from the tweets text

```
trump_tibble_extracted <- trump_tibble %>% unnest_tokens(word, text)
head(trump_tibble_extracted)
```

```
##           document word
## 1  01-19-2016 12:50:01  wow
## 1.1 01-19-2016 12:50:01  new
## 1.2 01-19-2016 12:50:01 polls
## 1.3 01-19-2016 12:50:01  just
## 1.4 01-19-2016 12:50:01  out
## 1.5 01-19-2016 12:50:01  have
```

(b) Create a dataframe consisting of the document-word counts

```
trump_word_counts <- trump_tibble_extracted %>%
  anti_join(stop_words) %>%
  count(document, word, sort = TRUE) %>%
  ungroup()
```

```
## Joining, by = "word"
```

(c) Create a document-term matrix using the `cast-dtm` function

```
trump_dtm <- trump_word_counts %>%
  cast_dtm(document, word, n)
trump_dtm
```

```
## <<DocumentTermMatrix (documents: 2725, terms: 5352)>>
## Non-/sparse entries: 23138/14561062
## Sparsity           : 100%
## Maximal term length: 38
## Weighting          : term frequency (tf)
```

## 2. LDA

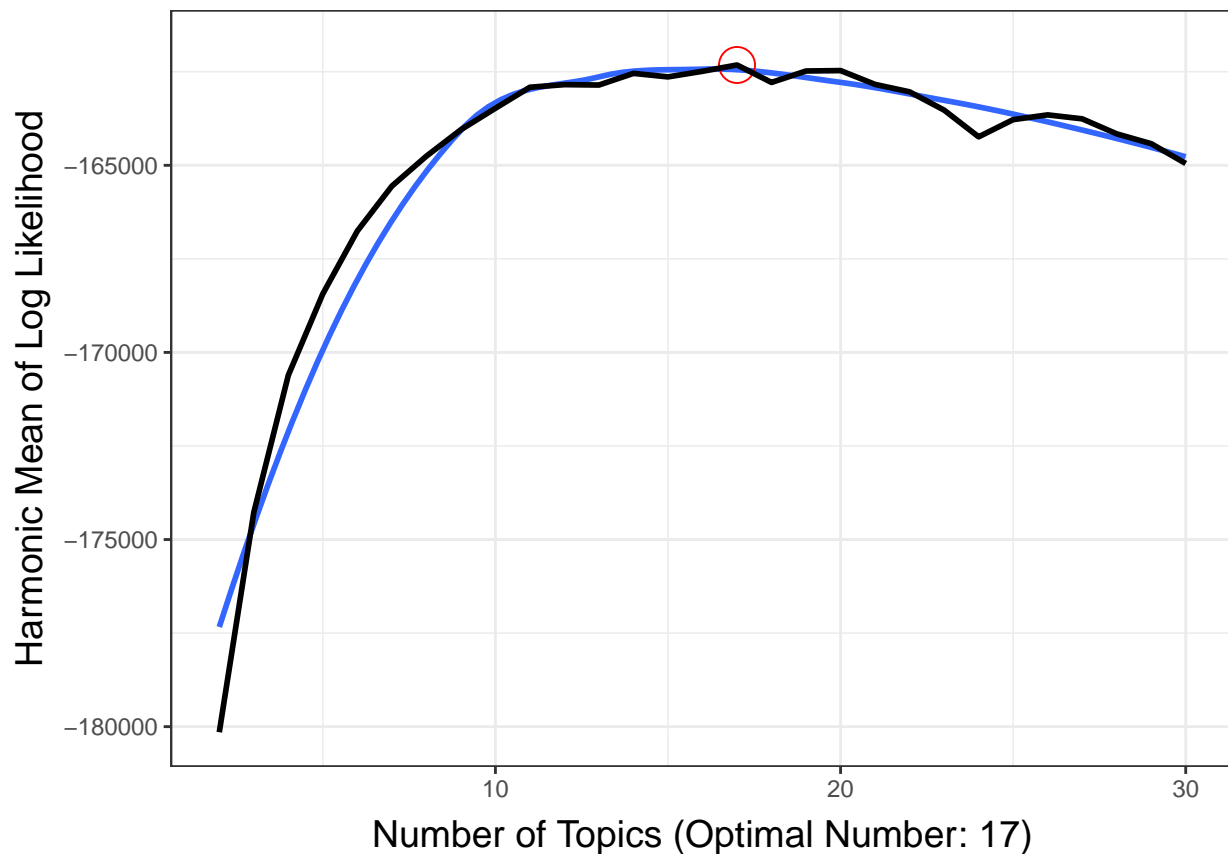
(a) Using the following control parameters, run the `optimal-k` function to search for the optimal number of topics. Be sure to set the “`max.k`” parameter equal to 30.

```
control <- list(burnin = 500, iter = 1000, keep = 100, seed = 46)
lda.opt.k <- optimal_k(trump_dtm, max.k=30, control=control, drop.seed=FALSE)
```

```
##
## Grab a cup of coffee this could take a while...
## 10 of 30 iterations (Current: 11:01:15; Elapsed: .3 mins)
## 20 of 30 iterations (Current: 11:01:45; Elapsed: .8 mins; Remaining: ~.8 mins)
## 30 of 30 iterations (Current: 11:02:26; Elapsed: 1.5 mins; Remaining: ~0 mins)
## Optimal number of topics = 17
```

(b) Plot the results of the `optimal-k` function. What does this plot suggest about the number of topics in the text?

```
plot(lda.opt.k)
```



#### ANSWER (2. LDA (b)):

The plot above suggests that there are 17 topics in the text.

- (c) Run LDA on the document-term matrix using the optimal value of  $k$ . Print out the top 10 words for each of the  $k$  topics. Comment on the results and their plausibility.

```
trump_lda = LDA(trump_dtm, k = as.numeric(lda.opt.k),
               method="Gibbs", control=control)
lda_inf = posterior(trump_lda)
topics.trump = topics(trump_lda, 1)
terms.trump = terms(trump_lda, 10)
print(terms.trump[,])
```

```
##      Topic 1  Topic 2  Topic 3  Topic 4  Topic 5
## [1,] "jobs"    "trump"   "enjoy"   "people" "president"
## [2,] "tax"     "honor"   "tonight" "american" "job"
## [3,] "cuts"    "donald"  "interviewed" "day" "money"
## [4,] "u.s"     "president" "7" "totally" "nice"
## [5,] "love"    "phony"   "00" "world" "rubio"
## [6,] "economy" "washington" "foxnews" "change" "poll"
## [7,] "stock"   "congress" "wonderful" "3" "meeting"
## [8,] "market"  "leaving"  "foxandfriends" "elizabeth" "families"
## [9,] "cut"     "national" "p.m" "warren" "law"
## [10,] "record" "flag"    "press" "incredible" "movement"
##      Topic 6  Topic 7  Topic 8  Topic 9
## [1,] "trump2016" "hillary" "republican" "watch"
## [2,] "makeamericagreatagain" "clinton" "house" "women"
## [3,] "rally" "crooked" "healthcare" "dollars"
```

```

## [4,] "arizona"          "bad"      "dems"      "trump2016"
## [5,] "votetrump"        "time"     "senate"    "record"
## [6,] "indiana"          "bernie"   "republicans" "time"
## [7,] "york"             "wow"      "obamacare" "russia"
## [8,] "virginia"         "run"      "bill"      "negative"
## [9,] "u.s"              "fbi"      "replace"   "ad"
## [10,] "hampshire"       "close"    "repeal"    "weak"
##      Topic 10 Topic 11 Topic 12 Topic 13 Topic 14 Topic 15
## [1,] "north"      "security" "america"   "cruz"     "join"     "vote"
## [2,] "korea"      "military" "day"       "ted"      "tomorrow" "election"
## [3,] "hard"       "u.s"      "speech"    "lyin"     "maga"     "florida"
## [4,] "people"     "democrats" "texas"     "campaign" "ohio"     "alabama"
## [5,] "support"    "border"   "forward"   "kasich"   "america"  "luther"
## [6,] "china"      "deal"     "happy"     "win"      "time"     "time"
## [7,] "carolina"   "trade"    "trump2016" "candidate" "watch"    "smart"
## [8,] "lost"       "wall"     "usa"       "jeb"      "night"    "strange"
## [9,] "south"      "sad"      "louisiana" "beat"     "live"     "8"
## [10,] "stop"      "mexico"   "massive"   "failed"   "fitn"     "russia"
##      Topic 16 Topic 17
## [1,] "news"      "country"
## [2,] "fake"      "people"
## [3,] "media"     "leaders"
## [4,] "cnn"       "remember"
## [5,] "story"     "united"
## [6,] "dishonest" "presidential"
## [7,] "failing"   "proud"
## [8,] "nytimes"   "spent"
## [9,] "lives"     "isis"
## [10,] "total"    "2017"

```

### ANSWER (2. LDA (c)):

The results are consistent with the assumption that each topic can be viewed as a different distribution of terms. The top 10 words vary by topic. These topics make sense as they reflect well what Trump focused on discussing since he took the political stage. For instance, topic 1 is concerned with economics including top words as "jobs", "tax", "economy" and "stock"; topic 3 reflects the fact that Trump is rumored to only get his news from Fox News and that Fox News is very Trump-friendly; topic 6 is concerned with his rally speeches including his famous campaign slogan "make America great again"; topic 7 reflects his constant bashing of his political opponent Hilary Clinton; topic 8 is more about politics including "republican", "house", and "repeal"; topic 10 is about Eastern Asia relationships including "north", "korea", "china" and "south"; and topic 16 is his "fake news" claim about the major media outlets.

## 3. Bayesian Logistic Regression

The first model we will fit to the contraceptives data is a varying-intercept logistic regression model, where the intercept varies by district.

Prior distribution:

$$\beta_{0j} \sim N(\mu_0, \sigma_0), \text{ with } \mu_0 \sim N(0, 100) \text{ and } \sigma_0 \sim \text{Exponential}(.1)$$

$$\beta_1 \sim N(0, 100), \beta_2 \sim N(0, 100), \beta_3 \sim N(0, 100)$$

Model for data:

$$Y_{ij} \sim \text{Bernoulli}(p_{ij})$$

$$\text{logit } p_{ij} = \beta_{0j} + \beta_1 * \text{urban} + \beta_2 * \text{living-children} + \beta_3 * \text{age-mean}$$

where  $Y_{ij}$  is 1 if woman  $i$  in district  $j$  uses contraception, and 0 otherwise, and where  $i = 1, \dots, N$  and  $j = 1, \dots, J$  ( $N$  is the number of observations in the data, and  $J$  is the number of districts). The above notation assumes  $N(\mu, \sigma)$  is a normal distribution with mean  $\mu$  and standard deviation  $\sigma$ . Also, the above notation assumes  $\text{Exponential}(\lambda)$  has mean  $1/\lambda$ . These are consistent with the parameterizations in Stan.

After you read the train and test data into R, the following code will help with formatting:

```
# read in data
train <- read.csv("data/dataset_2_train.txt")
test <- read.csv("data/dataset_2_test.txt")

# convert everything to numeric
for (i in 1:ncol(train)) {
  train[,i] <- as.numeric(as.character(train[,i]))
  test[,i] <- as.numeric(as.character(test[,i]))
}

# map district 61 to 54 (so that districts are in order)
train_bad_indices <- which(train$district == 61)
train[train_bad_indices, 1] <- 54
test_bad_indices <- which(test$district == 61)
test[test_bad_indices, 1] <- 54
```

- (a) To verify the procedure, simulate binary response data (using the `rbinom` function) assuming the following parameter values (and using the existing features and district information from the training data):

```
library(boot)
library(dbplyr)
mu_beta_0 = 2
sigma_beta_0 = 1
set.seed(123) # to ensure the next line is common to everyone
beta_0 = rnorm(n=60, mean=mu_beta_0, sd=sigma_beta_0)
beta_1 = 4
beta_2 = -3
beta_3 = -2

logit_pij = rep(NA, nrow(train))

# simulate data
for (i in 1:nrow(train)) {
  logit_pij[i] = beta_0[train[i,]$district] +
    beta_1 * train[i,]$urban +
    beta_2 * train[i,]$living.children +
    beta_3 * train[i,]$age_mean
}

pij = inv.logit(logit_pij)

y_simulated = rbinom(n=nrow(train), size=1, prob=pij)
```

- (b) Fit the varying-intercept model specified above to your simulated data

```
library(rstan)
library(ggplot2)
library(bayesplot)
```

```

theme_set(bayesplot::theme_default())

# create list
stan_list3b <- list()
stan_list3b$Y <- y_simulated
stan_list3b$N <- length(y_simulated) # number of obs
stan_list3b$J <- length(unique(train$district)) # number of district
stan_list3b$district <- train$district
stan_list3b$urban <- train$urban
stan_list3b$living_children <- train$living.children
stan_list3b$age_mean <- train$age_mean

```

```

# stan code
stan_code3b <- c("
data {
  // Number of observations
  int N;
  // Number of districts
  int J;
  // List of features, one for each observation
  int district[N];
  int<lower=0, upper=1> urban[N];
  int living_children[N];
  real age_mean[N];
  // Binary response (integer array)
  int<lower=0, upper=1> Y[N];
}

parameters {
  real mu_0;
  real<lower=0> sigma_0;
  real beta_0j[J];
  real beta_1;
  real beta_2;
  real beta_3;
}

model {
  // Prior
  mu_0 ~ normal(0,100);
  sigma_0 ~ exponential(0.1);
  beta_1 ~ normal(0,100);
  beta_2 ~ normal(0,100);
  beta_3 ~ normal(0,100);

  // J different beta_0j priors
  for (j in 1:J) {
    beta_0j[j] ~ normal(mu_0, sigma_0);
  }

  // Likelihood
  for (n in 1:N) {
    Y[n] ~ bernoulli_logit(beta_0j[district[n]] +

```

```

        beta_1*urban[n] + beta_2*living_children[n] + beta_3*age_mean[n]);
    }
}

generated quantities {
  int y_rep[N];          // Draws from posterior predictive dist

  for (n in 1:N) {
    y_rep[n] = bernoulli_rng(inv_logit(beta_0j[district[n]] +
        beta_1*urban[n] + beta_2*living_children[n] + beta_3*age_mean[n]));
  }
}

")

```

```

# fit the model
options(mc.cores = parallel::detectCores())
fit3b <- stan(model_code = stan_code3b,
  data = stan_list3b,
  iter = 2000,
  chains = 4,
  seed = 46,
  refresh = FALSE)

```

```

## In file included from filef92493e2f10.cpp:8:
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/BH/include/boost/math/
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/BH/include/boost/confi
## /Library/Frameworks/R.framework/Versions/3.4/Resources/library/BH/include/boost/config/compiler/clang.hpp:
## # define BOOST_NO_CXX11_RVALUE_REFERENCES
##      ^
## <command line>:6:9: note: previous definition is here
## #define BOOST_NO_CXX11_RVALUE_REFERENCES 1
##      ^
## In file included from filef92493e2f10.cpp:8:
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/RcppEigen/include/Eige
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/RcppEigen/include/Eige
## /Library/Frameworks/R.framework/Versions/3.4/Resources/library/RcppEigen/include/Eigen/src/Core/util/Reena
##      #pragma clang diagnostic pop
##      ^
## In file included from filef92493e2f10.cpp:8:
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st

```









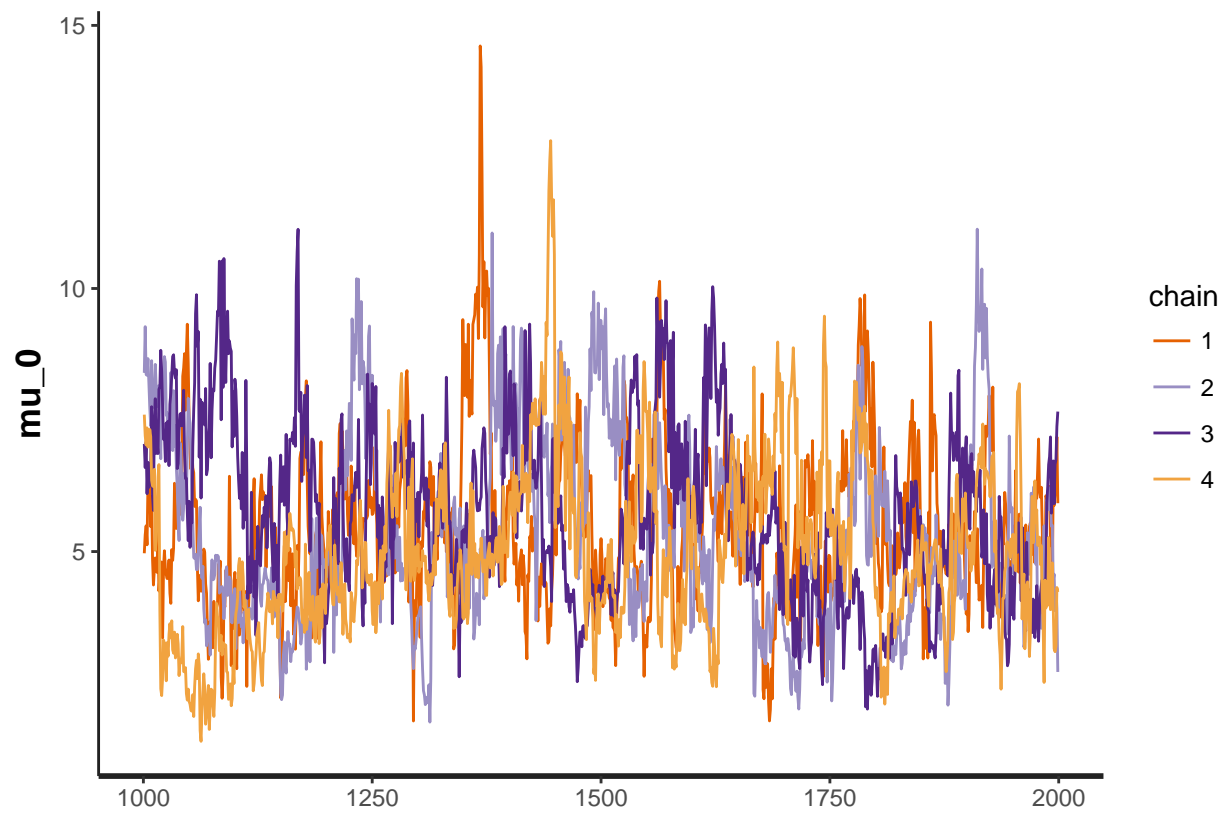
```

## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/RcppEigen/include/Eige
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/RcppEigen/include/Eige
## /Library/Frameworks/R.framework/Versions/3.4/Resources/library/RcppEigen/include/Eigen/src/Core/util/Reena
##      #pragma clang diagnostic pop
##      ^
## In file included from filef92493e2f10.cpp:629:
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/rstan/include/rstan/rs
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/rstan/include/rstan/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/sr
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/sr
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/sr
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/BH/include/boost/circu
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/BH/include/boost/circu
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/BH/include/boost/move/
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/BH/include/boost/move/
## /Library/Frameworks/R.framework/Versions/3.4/Resources/library/BH/include/boost/move/detail/iterator_trait
## BOOST_MOVE_STD_NS_BEG
## ^
## /Library/Frameworks/R.framework/Versions/3.4/Resources/library/BH/include/boost/move/detail/std_ns_begin.h
##      #define BOOST_MOVE_STD_NS_BEG _LIBCPP_BEGIN_NAMESPACE_STD
##      ^
## /Library/Developer/CommandLineTools/usr/include/c++/v1/__config:390:52: note: expanded from macro '_LIBCPP_F
## #define _LIBCPP_BEGIN_NAMESPACE_STD namespace std {inline namespace _LIBCPP_NAMESPACE {
##      ^
## In file included from filef92493e2f10.cpp:8:
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/sr
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/stan/math/rev/core/set_
##      static void set_zero_all_adjoints() {
##      ^
## In file included from filef92493e2f10.cpp:8:
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/sr
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/stan/math/rev/core/set_
##      static void set_zero_all_adjoints_nested() {
##      ^
## In file included from filef92493e2f10.cpp:8:
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/sr
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/stan/math/prim/mat/fun/
##      size_t fft_next_good_size(size_t N) {
##      ^
## 19 warnings generated.

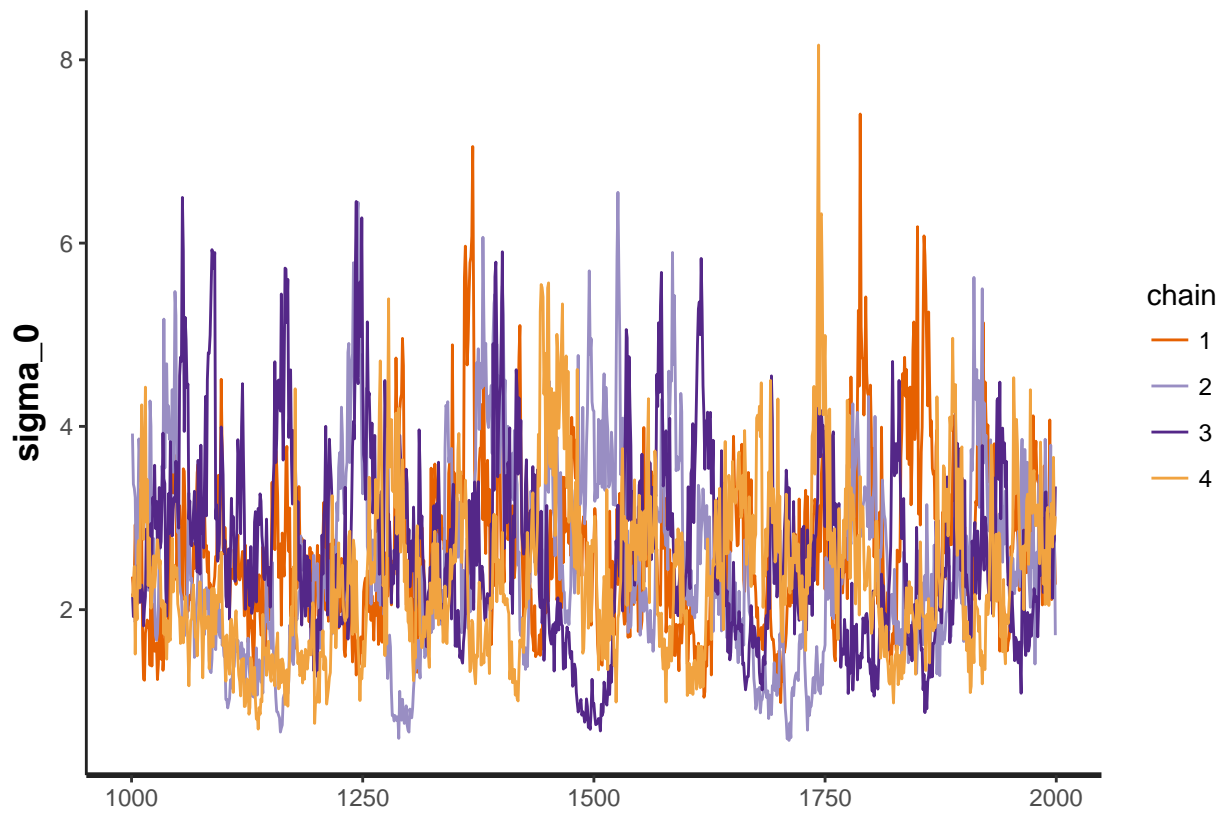
```

- (c) Plot the trace plots of the MCMC sampler for the parameters  $\mu_{\beta_0}, \sigma_{\beta_0}, \beta_1, \beta_2, \beta_3$ . Does it look like the samplers converged?

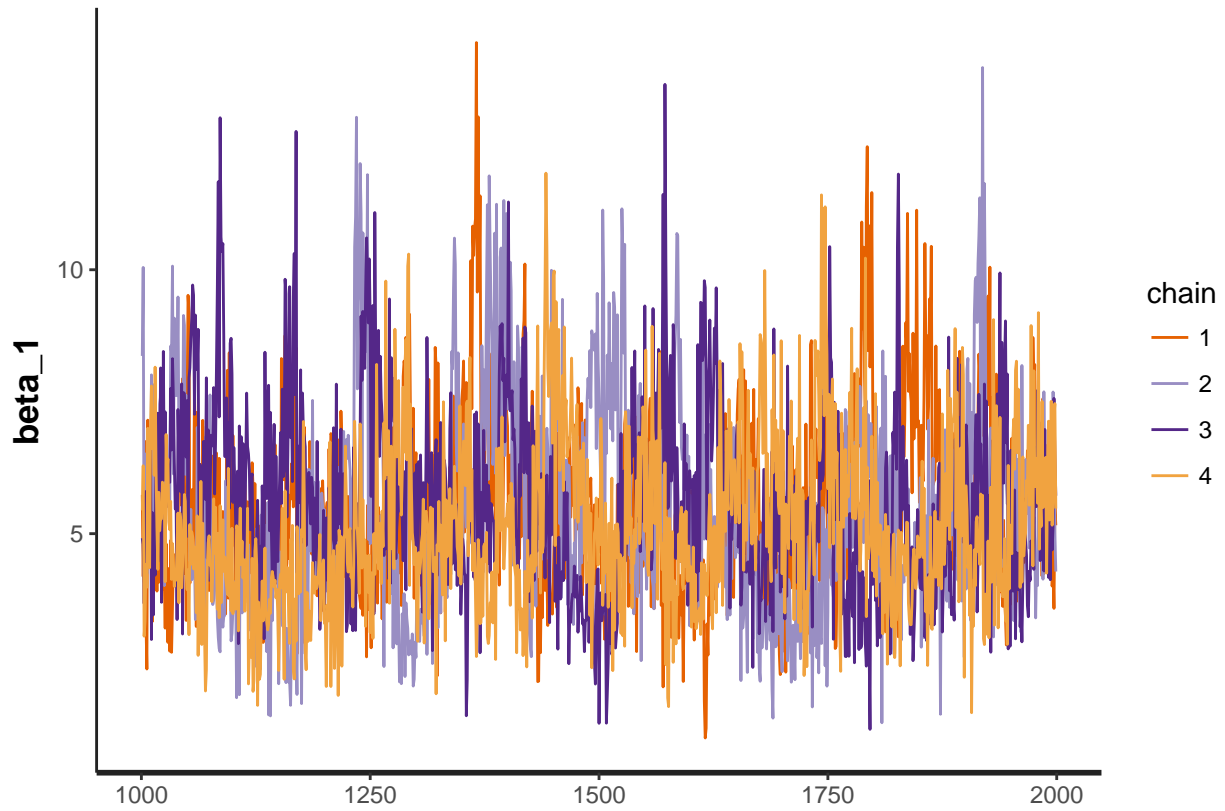
```
plot(fit3b, plotfun="trace", pars='mu_0')
```



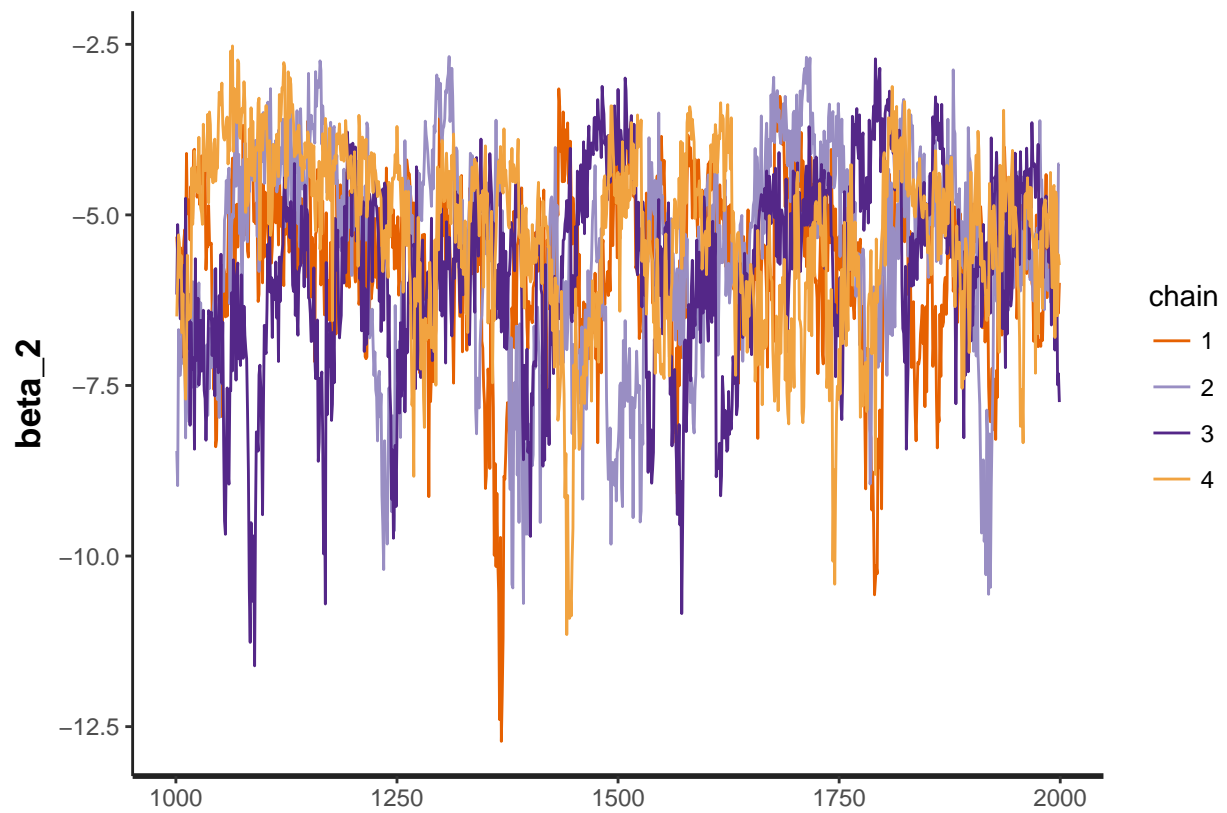
```
plot(fit3b, plotfun="trace", pars='sigma_0')
```



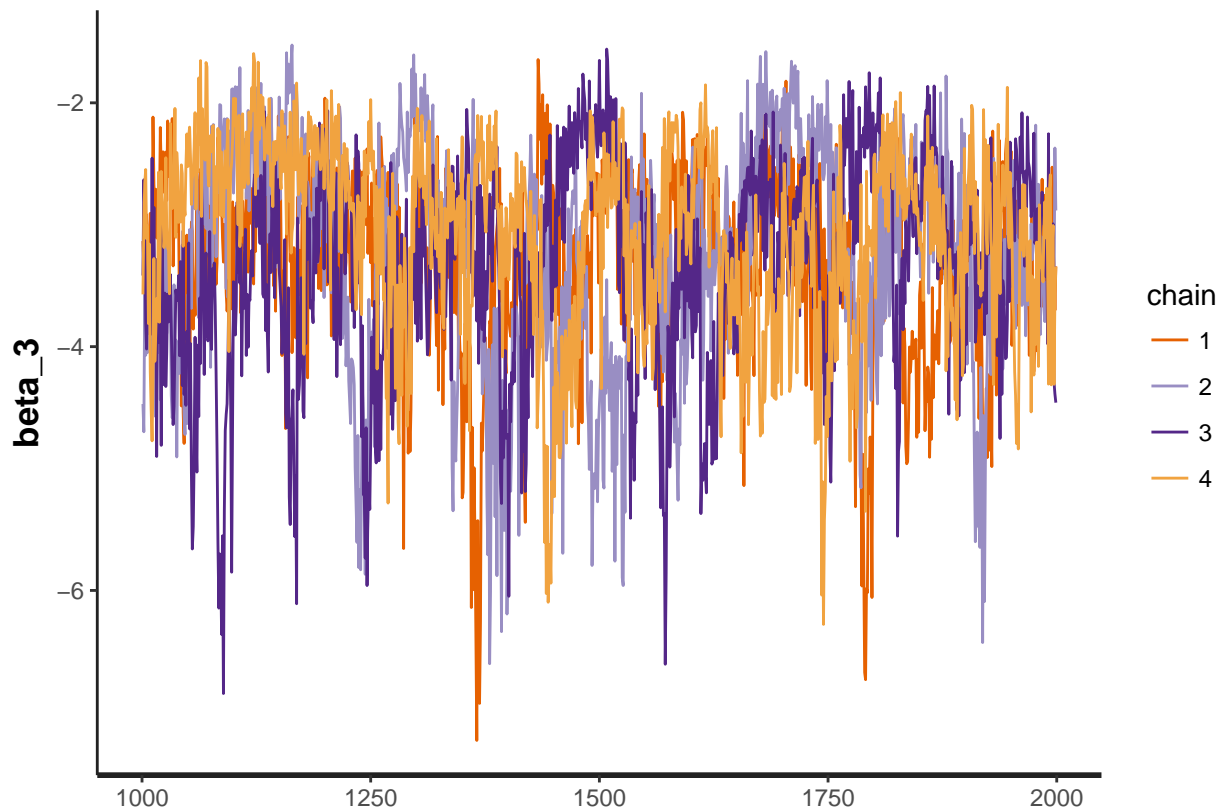
```
plot(fit3b, plotfun="trace", pars='beta_1')
```



```
plot(fit3b, plotfun="trace", pars='beta_2')
```



```
plot(fit3b, plotfun="trace", pars='beta_3')
```



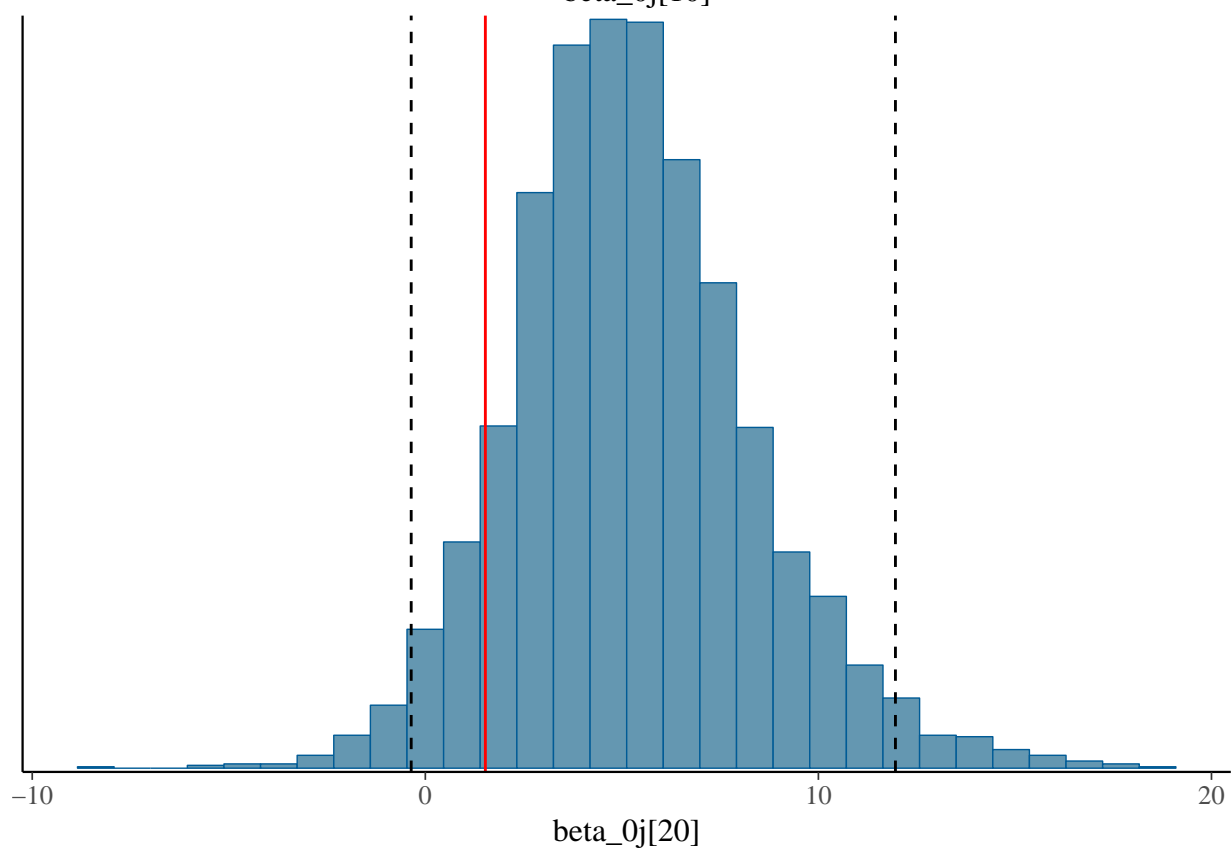
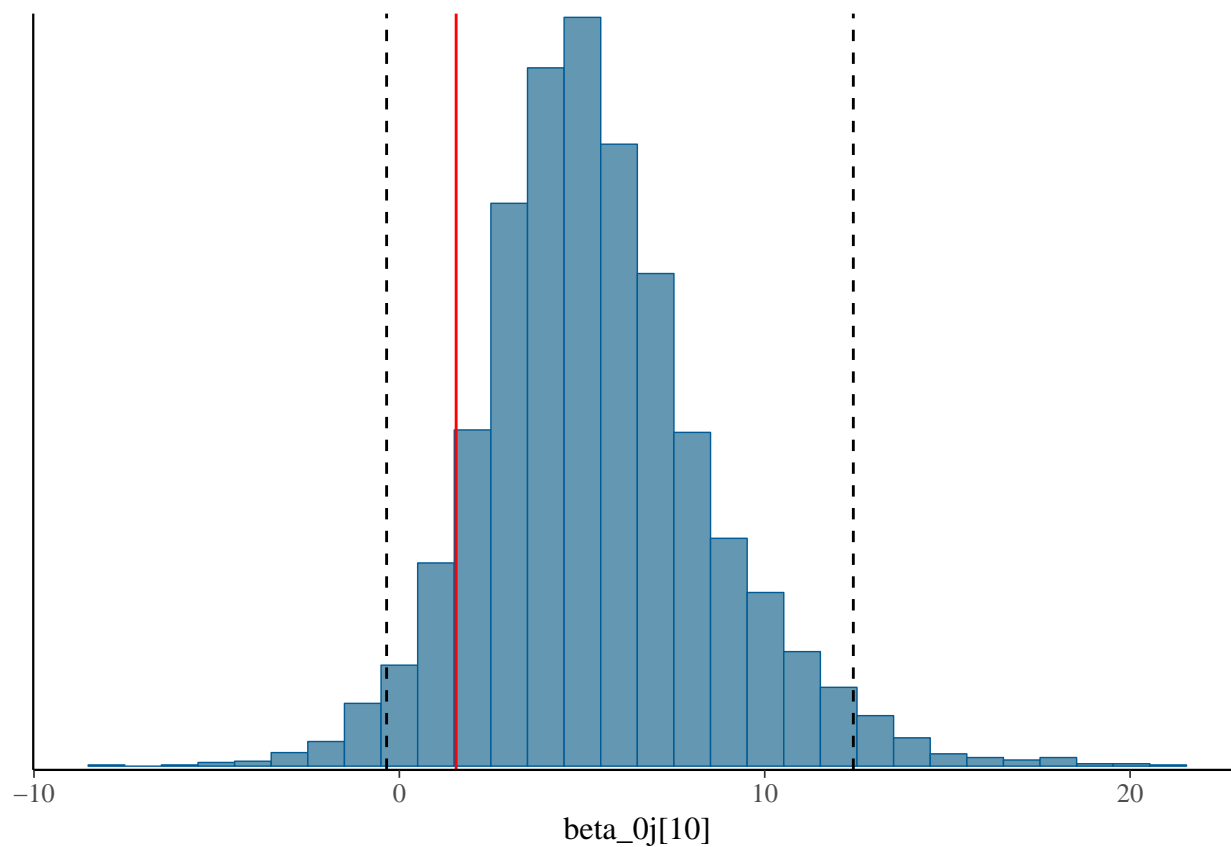
**ANSWER (3. Bayesian Logistic Regression (c)):**

The plot above suggests that the samples for the parameters (especially  $\mu_0$ ,  $\sigma_0$  and  $\beta_2$ ) did not converge very well.

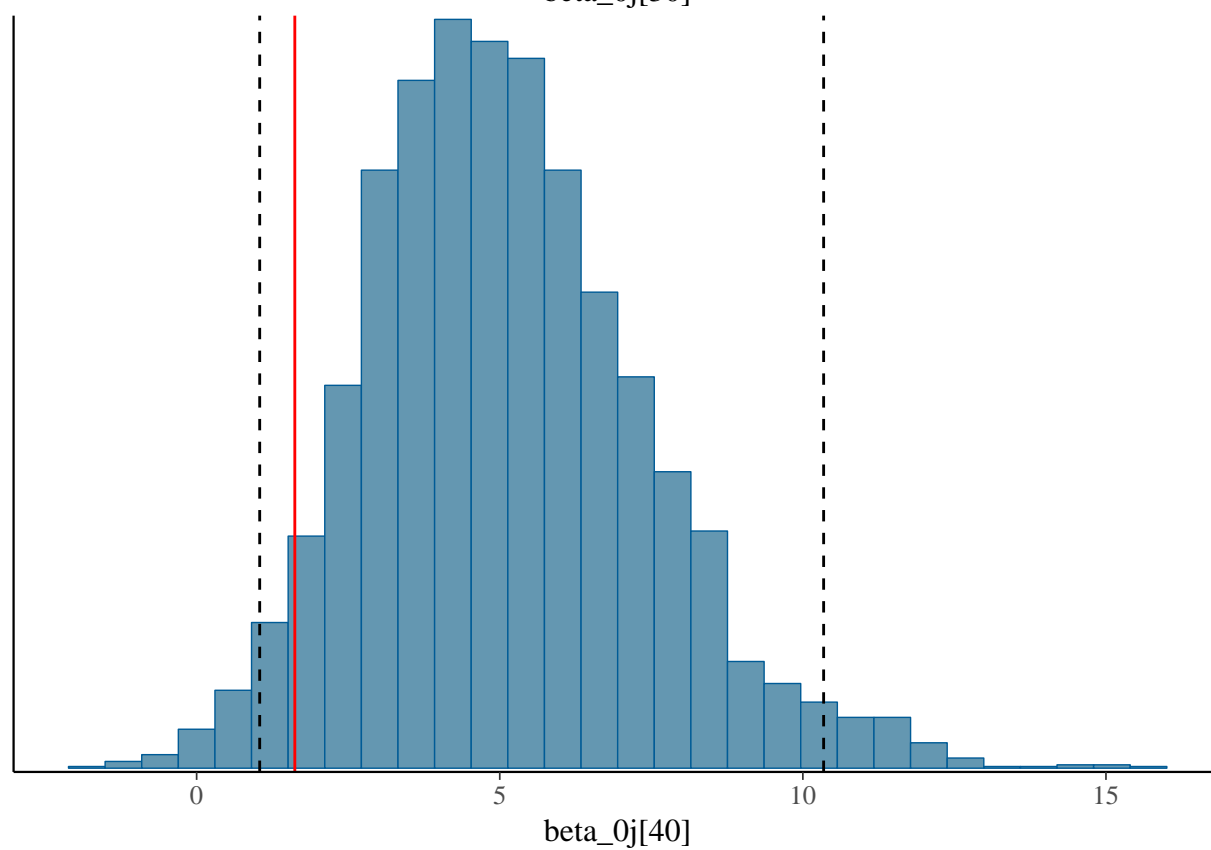
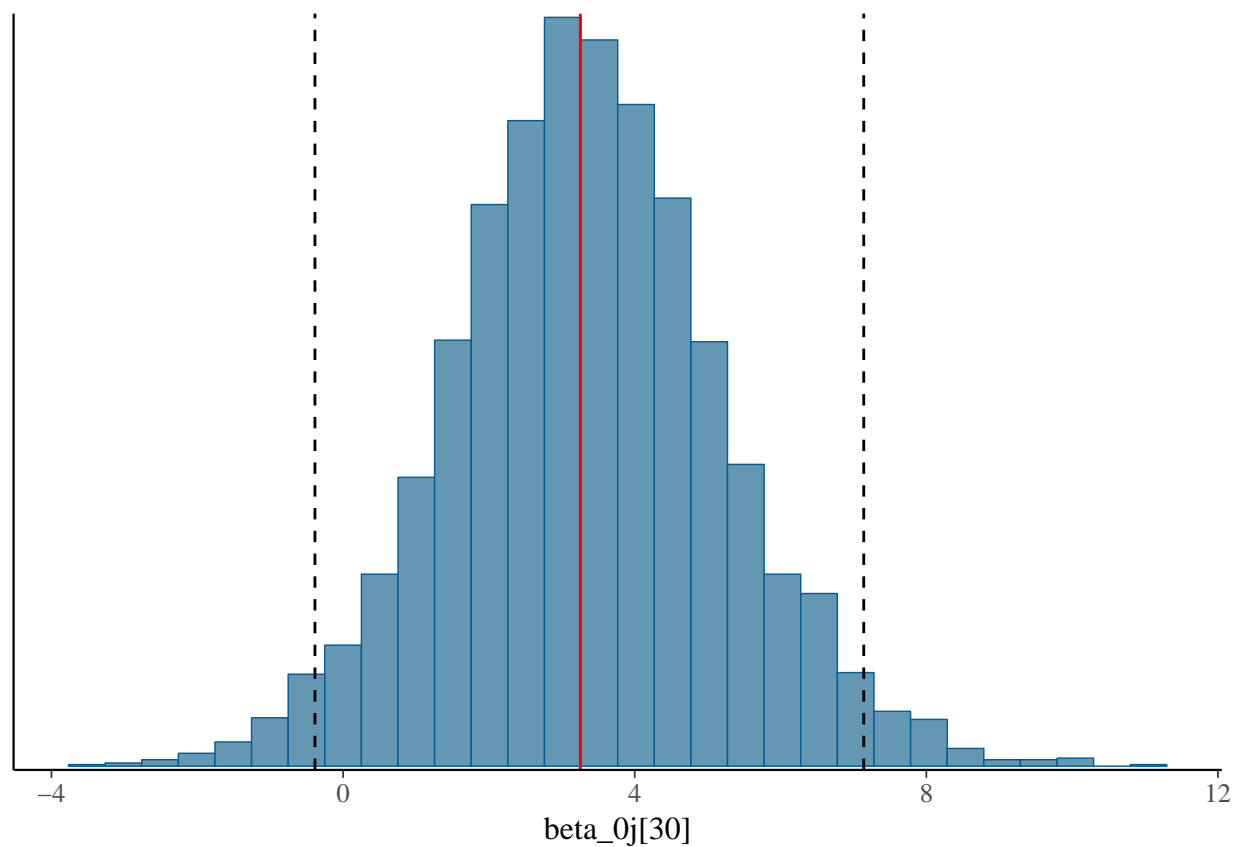
- (d) Plot histograms of the posterior distributions for the parameters  $\beta_{0,10}, \beta_{0,20}, \dots, \beta_{0,60}$ . Are the actual parameters that you generated contained within these posterior distributions?

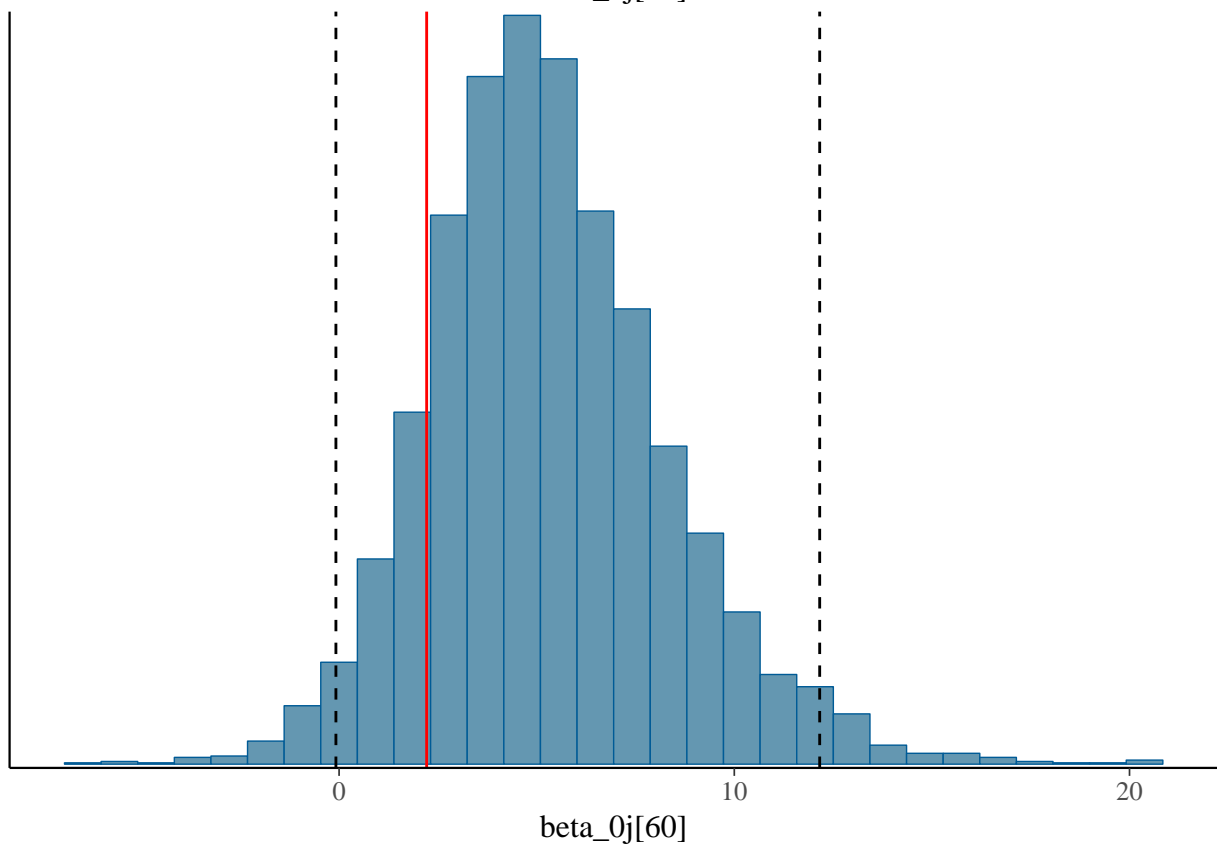
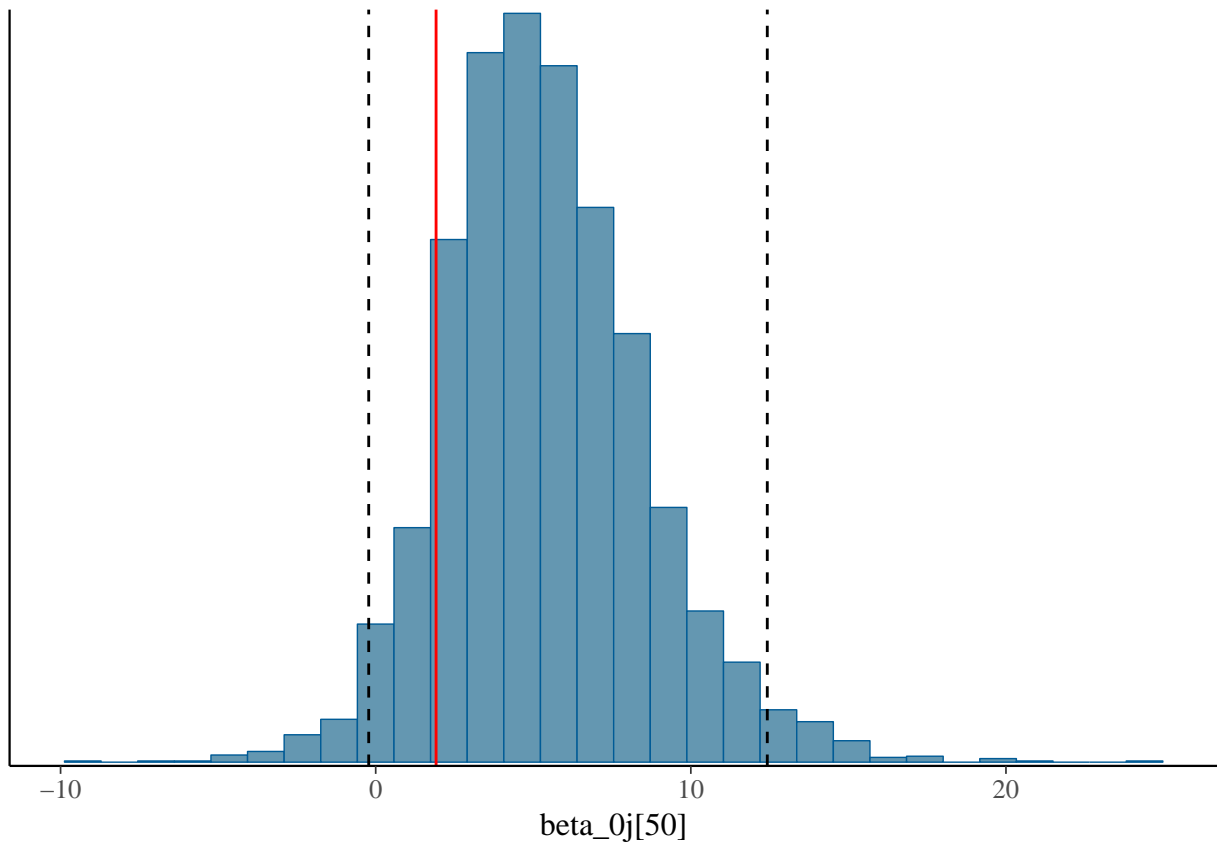
```
#color_scheme_set("brightblue") # check out bayesplot::color_scheme_set

for (i in seq(10,60,10)){
  beta_0j_draws <- as.matrix(fit3b, pars = paste0("beta_0j[", i, "]"))
  beta_0j_quantile <- quantile(beta_0j_draws, prob=c(0.025, 0.975))
  print(mcmc_hist(beta_0j_draws, prob = 0.95) +
    ggplot2::geom_vline(
      xintercept = beta_0[i], color="red"
    ) +
    ggplot2::geom_vline(
      xintercept = beta_0j_quantile, linetype="dashed"
    ))
}
```









ANSWER (3. Bayesian Logistic Regression (d)):

The plots above show the histograms of the posterior distributions for  $\beta_{0,10}, \beta_{0,20}, \dots, \beta_{0,60}$  with their 95% central interval shown in dotted black lines and the actual generated  $\beta_{0,10}, \beta_{0,20}, \dots, \beta_{0,60}$  shown in solid red line. Based on these plots, the actual generated parameters are all contained within the 95% central posterior intervals.

We now fit our model to the actual data.

- (e) Fit the varying-intercept model to the real train data. Make sure to set a seed at 46 within the Stan function, to ensure that you will get the same results if you fit your model correctly.

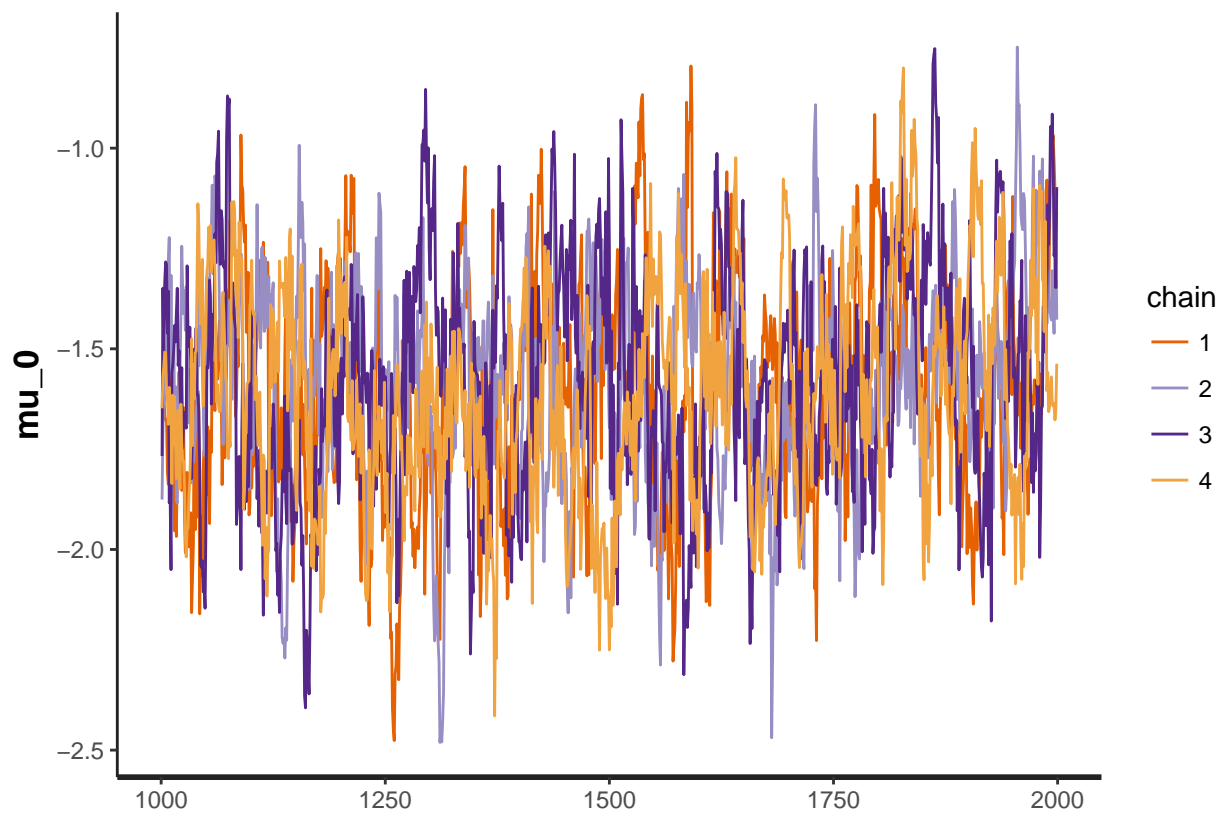
```
# create list
stan_list3e <- list()
stan_list3e$Y <- train$contraceptive_use
stan_list3e$N <- nrow(train) # number of obs
stan_list3e$J <- length(unique(train$district)) # number of district
stan_list3e$district <- train$district
stan_list3e$urban <- train$urban
stan_list3e$living_children <- train$living.children
stan_list3e$age_mean <- train$age_mean

stan_list3e$N_test <- nrow(test) # number of obs
stan_list3e$J_test <- length(unique(test$district)) # number of district
stan_list3e$district_test <- test$district
stan_list3e$urban_test <- test$urban
stan_list3e$living_children_test <- test$living.children
stan_list3e$age_mean_test <- test$age_mean

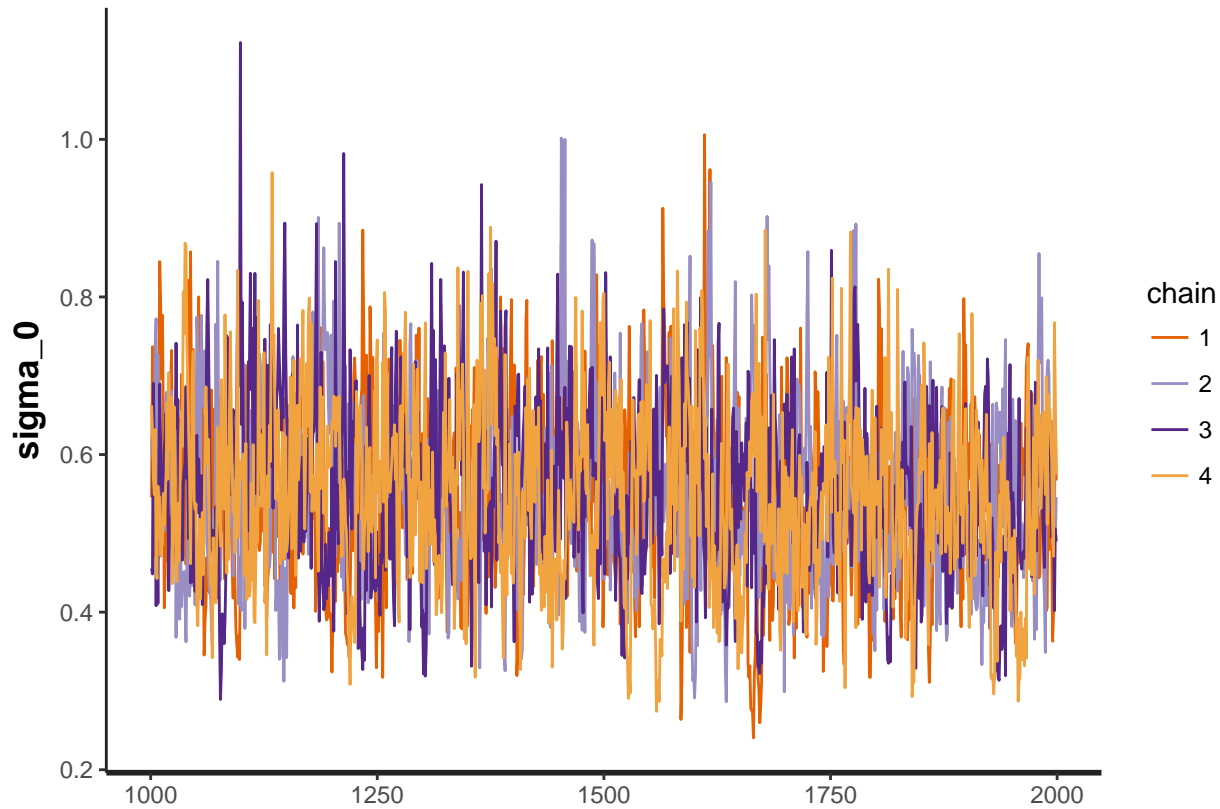
# fit the model
options(mc.cores = parallel::detectCores())
fit3e <- stan(model_code = stan_code3b,
  data = stan_list3e,
  iter = 2000,
  chains = 4,
  seed = 46,
  refresh = FALSE)
```

- (f) Check the convergence by examining the trace plots, as you did with the simulated data.

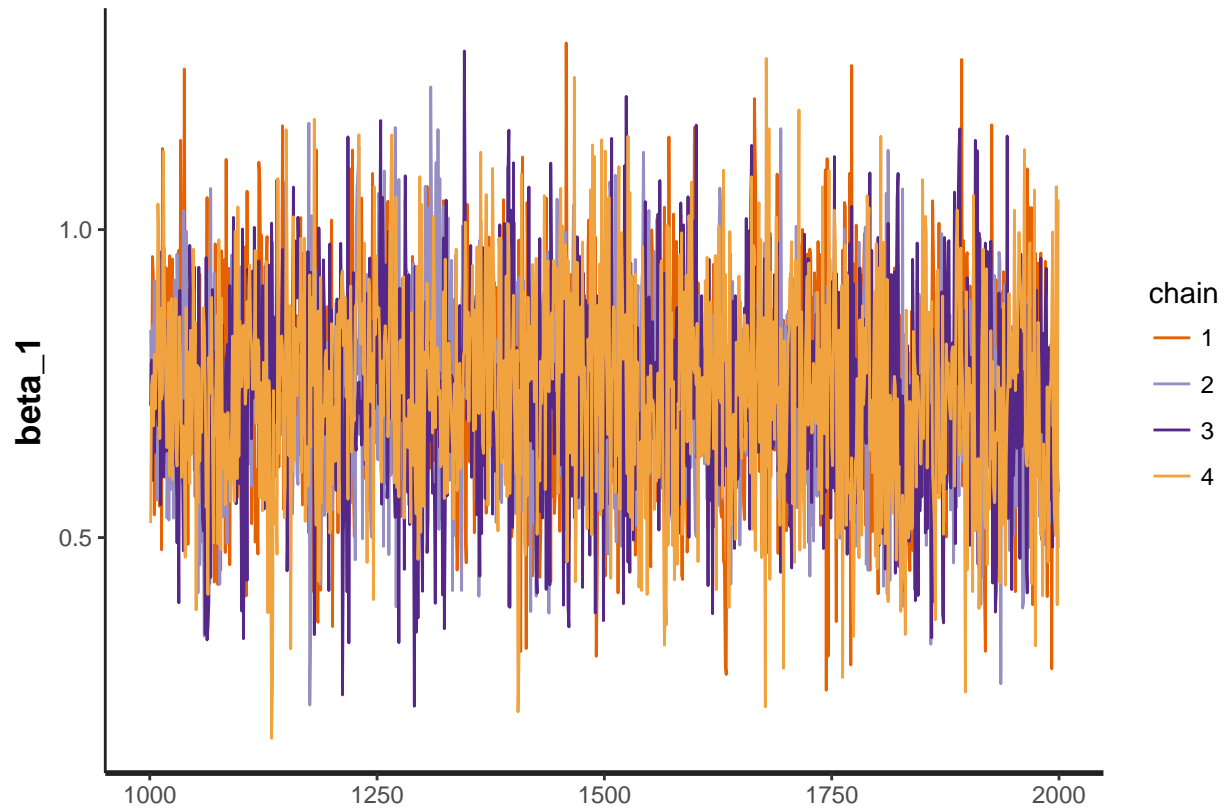
```
plot(fit3e, plotfun="trace", pars='mu_0')
```



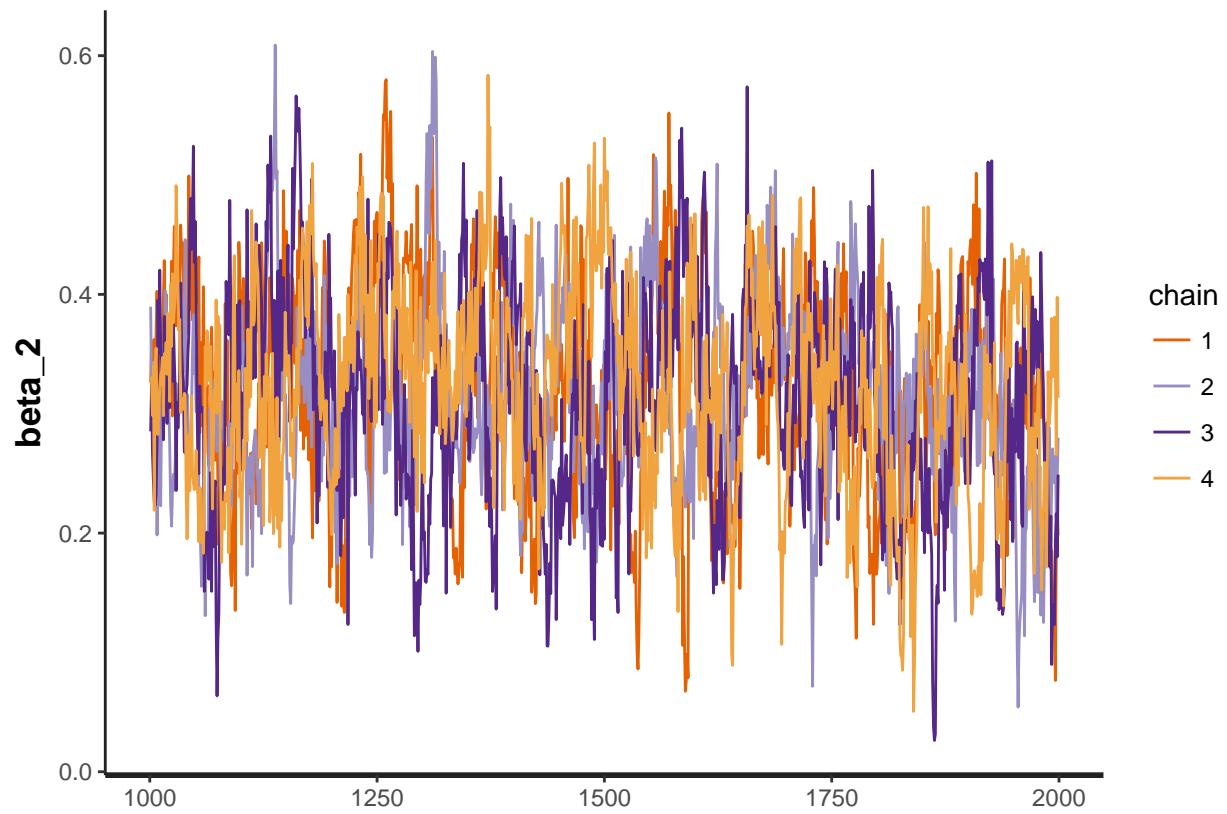
```
plot(fit3e, plotfun="trace", pars='sigma_0')
```



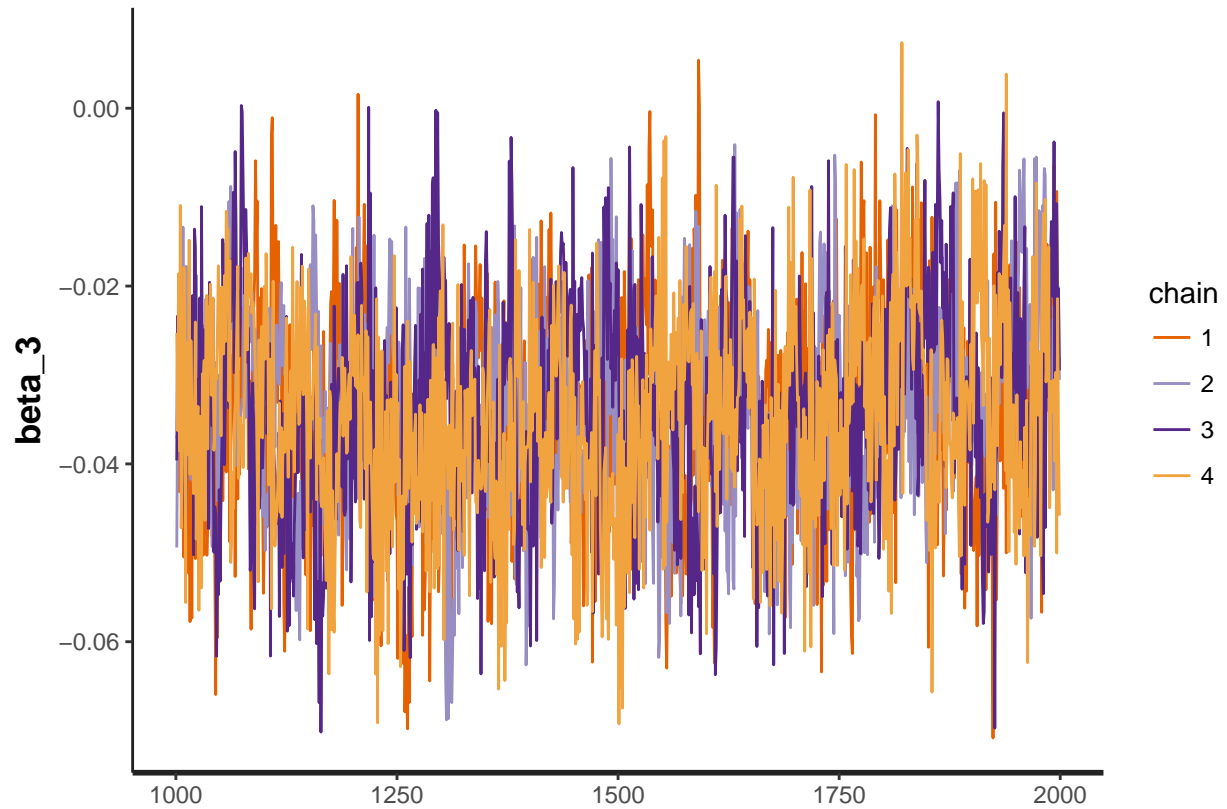
```
plot(fit3e, plotfun="trace", pars='beta_1')
```



```
plot(fit3e, plotfun="trace", pars='beta_2')
```



```
plot(fit3e, plotfun="trace", pars='beta_3')
```



**ANSWER (3. Bayesian Logistic Regression (f)):**

The plot above suggests that the samples for the parameters converged.

- (g) Based on the posterior means, women belonging to which district are most likely to use contraceptives?  
Women belonging to which district are least likely to use contraceptives?

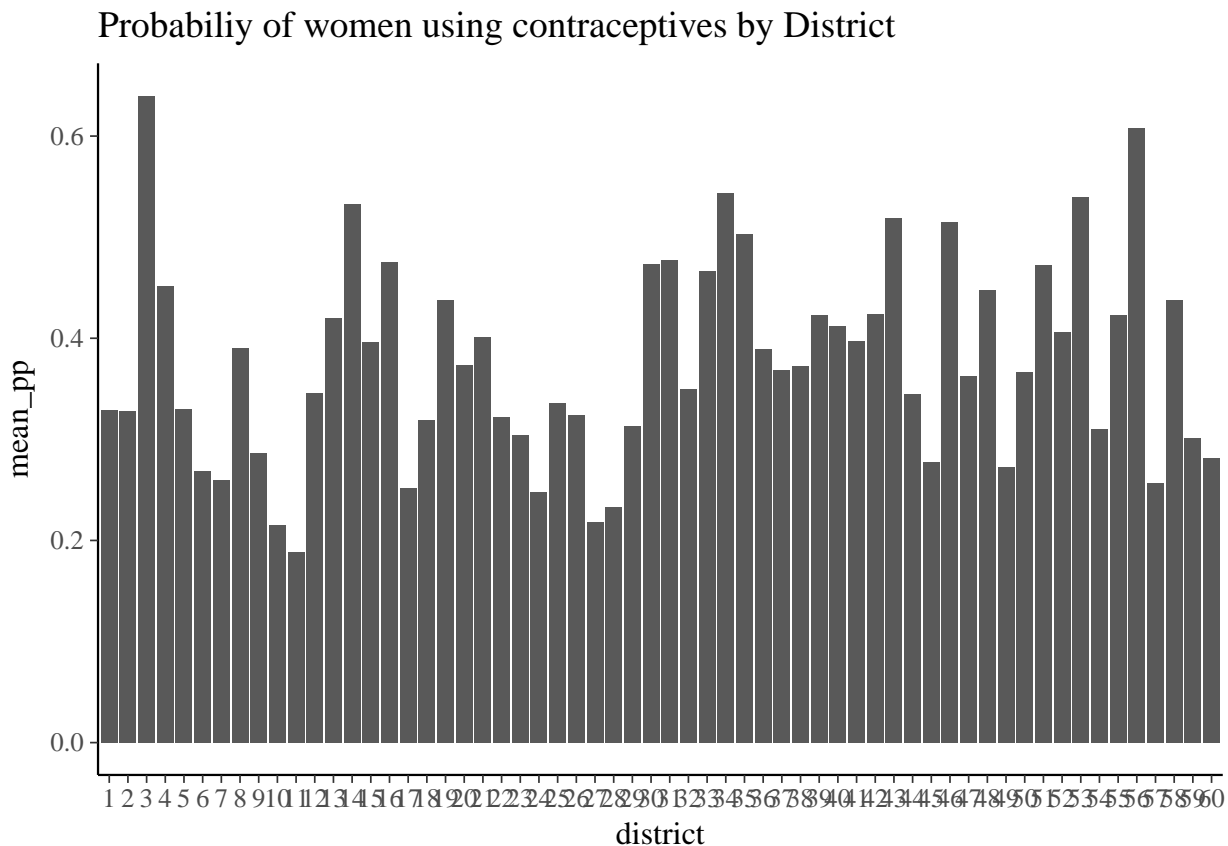
```
# mean of y_rep3e from 4000 iterations
y_rep3e_mean <- colMeans(as.matrix(fit3e, pars="y_rep"))

# make a data frame combining district information
df_yrep3e <- data.frame(cbind(district = train$district, y_rep3e_mean = y_rep3e_mean))

# compute mean of posterior predictive by district
mean_yrep3e_by_district <- df_yrep3e %>% group_by(district = factor(district)) %>%
  summarise(mean_pp = mean(y_rep3e_mean))

# plot average posterior predictive by district
mean_yrep3e_by_district.plot <- ggplot(data=mean_yrep3e_by_district, aes(x=district, y=mean_pp)) +
  geom_bar(stat="identity") + ggtitle("Probabiliy of women using contraceptives by District")

mean_yrep3e_by_district.plot
```



```
# get the district with maximum and minimum average probability
max_prob_district = mean_yrep3e_by_district$district[mean_yrep3e_by_district$mean_pp ==
  max(mean_yrep3e_by_district$mean_pp)]
print(paste0("most likely district: ", max_prob_district))
```

```
## [1] "most likely district: 3"
```

```

min_prob_district = mean_yrep3e_by_district$district[mean_yrep3e_by_district$mean_pp ==
                                                    min(mean_yrep3e_by_district$mean_pp)]
min_prob_district

## [1] 11
## 60 Levels: 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 ... 60
print(paste0("least likely district: ", min_prob_district))

## [1] "least likely district: 11"

```

### ANSWER (3. Bayesian Logistic Regression (g)):

Based on the R code and plot above, women belonging to 3 are most likely to use contraceptives while women belonging to 11 are least likely to use contraceptives.

- (h) What are the posterior means of  $\mu_{\beta_0}$  and  $\sigma_{\beta_0}$ ? Do these values offer any evidence in support of or against the varying-intercept model?

```

beta_0_summary <- summary(fit3e, pars=c("mu_0", "sigma_0"), probs=c(0.05, 0.95))$summary
mu_beta_0_mean = beta_0_summary[, "mean"][1]
sigma_beta_0_mean = beta_0_summary[, "mean"][2]

print(paste0('posterior mean of mu_beta_0 is: ', mu_beta_0_mean))

## [1] "posterior mean of mu_beta_0 is: -1.59388965418123"
print(paste0('posterior mean of sigma_beta_0 is: ', sigma_beta_0_mean))

```

```

## [1] "posterior mean of sigma_beta_0 is: 0.550823270039476"
mu_beta_0_draws_real <- as.matrix(fit3e, pars = "mu_0")
sigma_beta_0_draws_real <- as.matrix(fit3e, pars = "sigma_0")

```

### ANSWER (3. Bayesian Logistic Regression (h)):

The posterior mean of  $\mu_{\beta_0} = -1.5938897$ . The posterior mean of  $\sigma_{\beta_0} = 0.5508233$ . Therefore, the 95% confidence interval of  $\beta_0$  is  $\mu_{\beta_0} \pm 2 * \sigma_{\beta_0} = (-0.4922431, -2.6955362)$ . Since the 95% confidence interval does not capture 0, we have enough evidence in support of the varying-intercept model.

## 4. Varying-Coefficients Model

In the next model we will fit to the contraceptives data is a varying-coefficients logistic regression model, where the coefficient on living-children varies by district:

Prior distribution:

$$\beta_{0j} \sim N(\mu_0, \sigma_0), \text{ with } \mu_0 \sim N(0, 100) \text{ and } \sigma_0 \sim \text{Exponential}(0.1)$$

$$\beta_{1j} \sim N(0, \sigma_1), \text{ with } \sigma_1 \sim \text{Exponential}(0.1)$$

$$\beta_{2j} \sim N(0, \sigma_2), \text{ with } \sigma_2 \sim \text{Exponential}(0.1)$$

$$\beta_{3j} \sim N(0, \sigma_3), \text{ with } \sigma_3 \sim \text{Exponential}(0.1)$$

Model for data:

$$Y_{ij} \sim \text{Bernoulli}(p_{ij})$$

$$\text{logit } p_{ij} = \beta_{0j} + \beta_{1j}\text{urban} + \beta_{2j}\text{age-mean} + \beta_{3j}\text{living-children}$$

where  $i = 1, \dots, N$  and  $j = 1, \dots, J$  ( $N$  is the number of observations in the data, and  $J$  is the number of districts).



- (a) Fit the model to the real data. For each of the three coefficients to the predictors, plot vertical segments corresponding to the 95% central posterior intervals for the coefficient within each district. Thus you should have 54 parallel segments on each graph. If the segments are overlapping on the vertical scale, then the model fit suggests that the coefficient does not differ by district. What do you conclude from these graphs?

```
# stan code for problem 4
stan_code4a <- c("
data {
  // Number of observations
  int N;
  // Number of districts
  int J;
  // List of features, one for each observation
  int district[N];
  int<lower=0, upper=1> urban[N];
  int living_children[N];
  real age_mean[N];
  // Binary response (integer array)
  int<lower=0, upper=1> Y[N];

  int N_test;
  // Number of districts
  int J_test;
  // List of features, one for each observation
  int district_test[N_test];
  int<lower=0, upper=1> urban_test[N_test];
  int living_children_test[N_test];
  real age_mean_test[N_test];
}

parameters {
  real mu_0;
  real<lower=0> sigma_0;
  real beta_0j[J]; // bias term vary by J districts
  real beta_1j[J];
  real<lower=0> sigma_1;
  real beta_2j[J];
  real<lower=0> sigma_2;
  real beta_3j[J]; // living children vary by J districts
  real<lower=0> sigma_3;
}

model {
  // Prior
  mu_0 ~ normal(0,100);
  sigma_0 ~ exponential(0.1);
  sigma_1 ~ exponential(0.1);
  sigma_2 ~ exponential(0.1);
  sigma_3 ~ exponential(0.1);

  // J different beta_0j priors
  for (j in 1:J) {
    beta_0j[j] ~ normal(mu_0, sigma_0);
```

```

    beta_1j[j] ~ normal(0,sigma_1);
    beta_2j[j] ~ normal(0,sigma_2);
    beta_3j[j] ~ normal(0,sigma_3);
  }

  // Likelihood
  for (n in 1:N) {
    Y[n] ~ bernoulli_logit(beta_0j[district[n]] + beta_1j[district[n]]*urban[n] +
      beta_2j[district[n]]*age_mean[n] + beta_3j[district[n]]*living_children[n]);
  }
}

generated quantities {
  int y_rep[N_test];          // Draws from posterior predictive dist

  for (n in 1:N_test) {
    y_rep[n] = bernoulli_rng(inv_logit(beta_0j[district_test[n]] +
      beta_1j[district_test[n]]*urban_test[n] +
      beta_2j[district_test[n]]*age_mean_test[n] +
      beta_3j[district_test[n]]*living_children_test[n]));
  }
}

")

```

```

# fit the model
options(mc.cores = parallel::detectCores())
fit4a <- stan(model_code = stan_code4a,
  data = stan_list3e,
  iter = 2000,
  chains = 4,
  seed = 46,
  refresh = FALSE)

```

```

## In file included from filef9245f89974c.cpp:8:
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/BH/include/boost/math/
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/BH/include/boost/confi
## /Library/Frameworks/R.framework/Versions/3.4/Resources/library/BH/include/boost/config/compiler/clang.hpp:
## # define BOOST_NO_CXX11_RVALUE_REFERENCES
##      ^
## <command line>:6:9: note: previous definition is here
## #define BOOST_NO_CXX11_RVALUE_REFERENCES 1
##      ^
## In file included from filef9245f89974c.cpp:8:
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/sr
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st

```







```

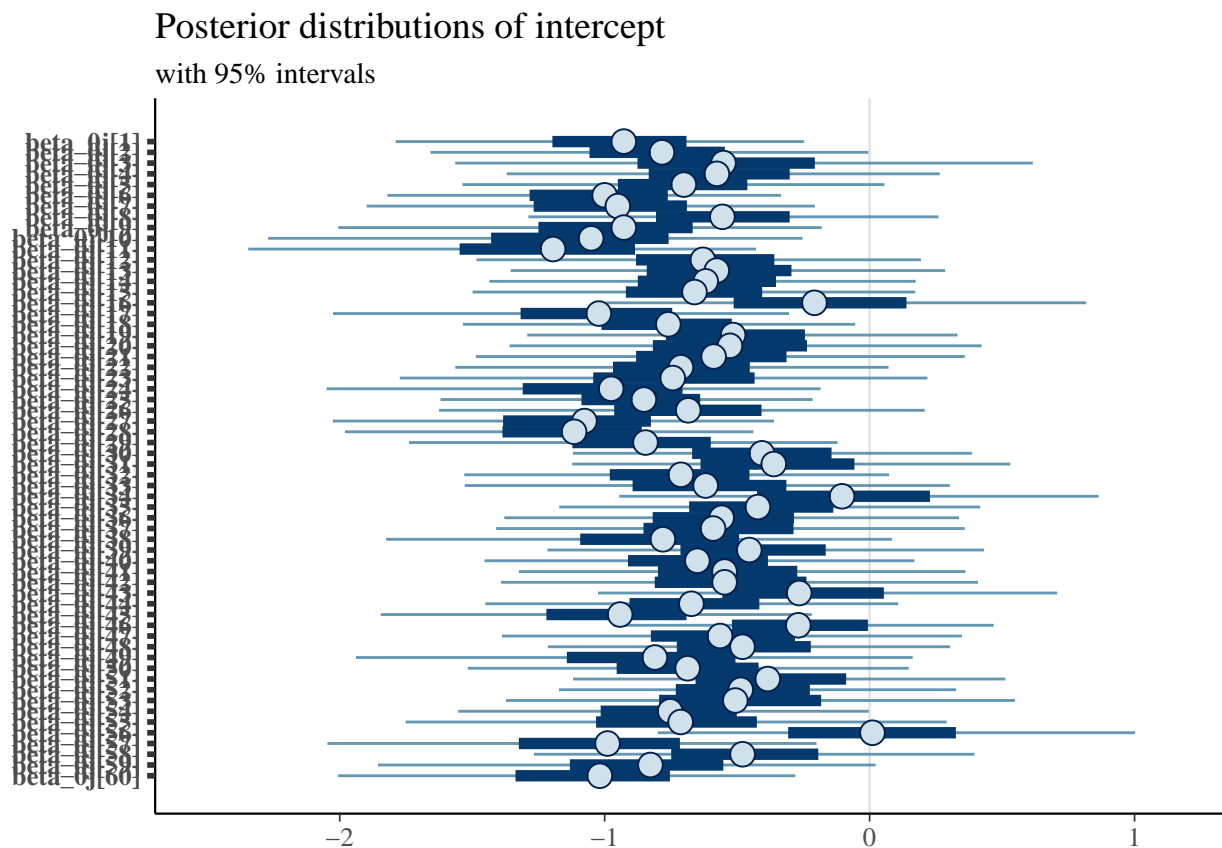
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/RcppEigen/include/Eige
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/RcppEigen/include/Eige
## /Library/Frameworks/R.framework/Versions/3.4/Resources/library/RcppEigen/include/Eigen/src/Core/util/Reena
##      #pragma clang diagnostic pop
##      ^
## In file included from filef9245f89974c.cpp:8:
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/sr
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/RcppEigen/include/Eige
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/RcppEigen/include/Eige
## /Library/Frameworks/R.framework/Versions/3.4/Resources/library/RcppEigen/include/Eigen/src/Core/util/Reena
##      #pragma clang diagnostic pop
##      ^
## In file included from filef9245f89974c.cpp:860:
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/rstan/include/rstan/rs
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/rstan/include/rstan/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/sr
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/sr
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/sr
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/BH/include/boost/circu
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/BH/include/boost/circu
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/BH/include/boost/move/
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/BH/include/boost/move/
## /Library/Frameworks/R.framework/Versions/3.4/Resources/library/BH/include/boost/move/detail/iterator_trait
## BOOST_MOVE_STD_NS_BEG
## ^
## /Library/Frameworks/R.framework/Versions/3.4/Resources/library/BH/include/boost/move/detail/std_ns_begin.h
##      #define BOOST_MOVE_STD_NS_BEG _LIBCPP_BEGIN_NAMESPACE_STD
##      ^
## /Library/Developer/CommandLineTools/usr/include/c++/v1/__config:390:52: note: expanded from macro '_LIBCPP_F
## #define _LIBCPP_BEGIN_NAMESPACE_STD namespace std {inline namespace _LIBCPP_NAMESPACE {
##      ^
## In file included from filef9245f89974c.cpp:8:
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/sr
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/stan/math/rev/core/set_
##      static void set_zero_all_adjoints() {
##      ^
## In file included from filef9245f89974c.cpp:8:
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/sr
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/stan/math/rev/core/set_
##      static void set_zero_all_adjoints_nested() {
##      ^
## In file included from filef9245f89974c.cpp:8:
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/sr

```

```
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/stan/math/prim/mat/fun/
##      size_t fft_next_good_size(size_t N) {
##      ^
## 19 warnings generated.

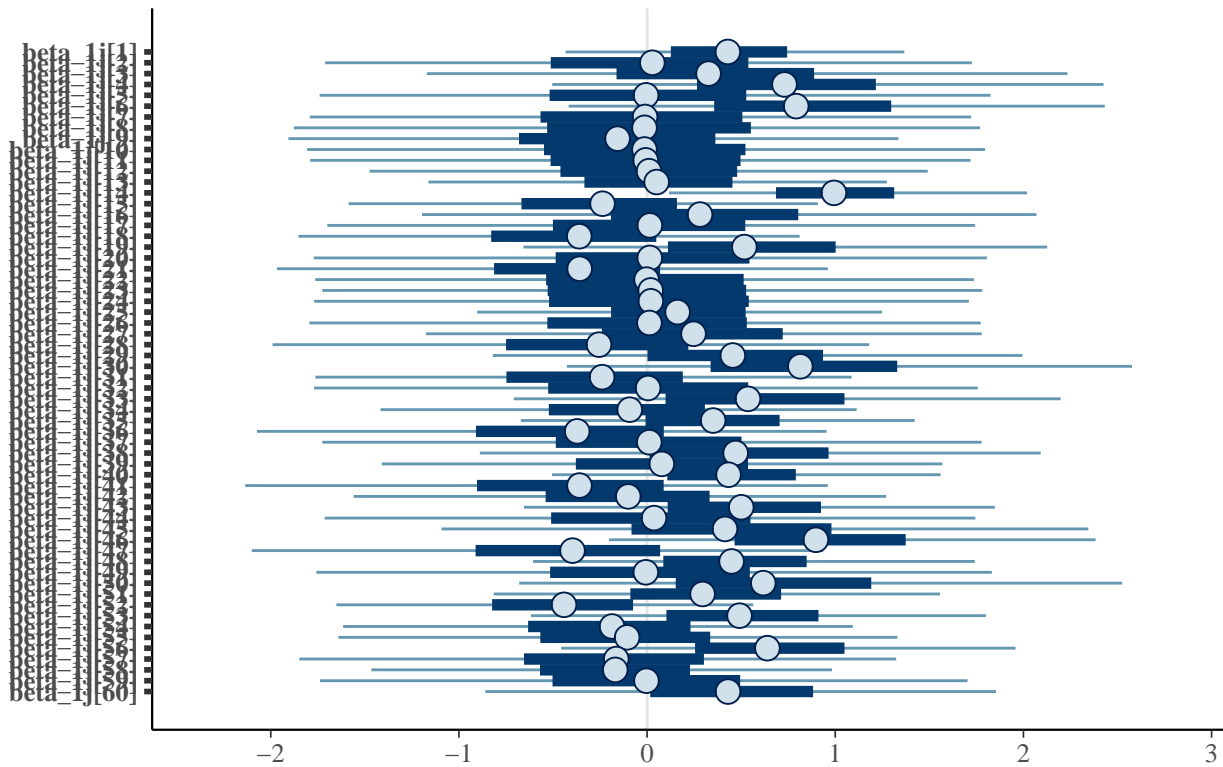
fit_posterior <- as.array(fit4a)
beta_0s <- sprintf("beta_0j[%d]", 1:60)
beta_1s <- sprintf("beta_1j[%d]", 1:60)
beta_2s <- sprintf("beta_2j[%d]", 1:60)
beta_3s <- sprintf("beta_3j[%d]", 1:60)

mcmc_intervals(fit_posterior, pars=beta_0s, prob_out=0.95) +
  ggplot2::labs(
    title = "Posterior distributions of intercept",
    subtitle = "with 95% intervals"
  )
)
```



```
mcmc_intervals(fit_posterior, pars=beta_1s, prob_out=0.95) +
  ggplot2::labs(
    title = "Posterior distributions of coef for urban",
    subtitle = "with 95% intervals"
  )
)
```

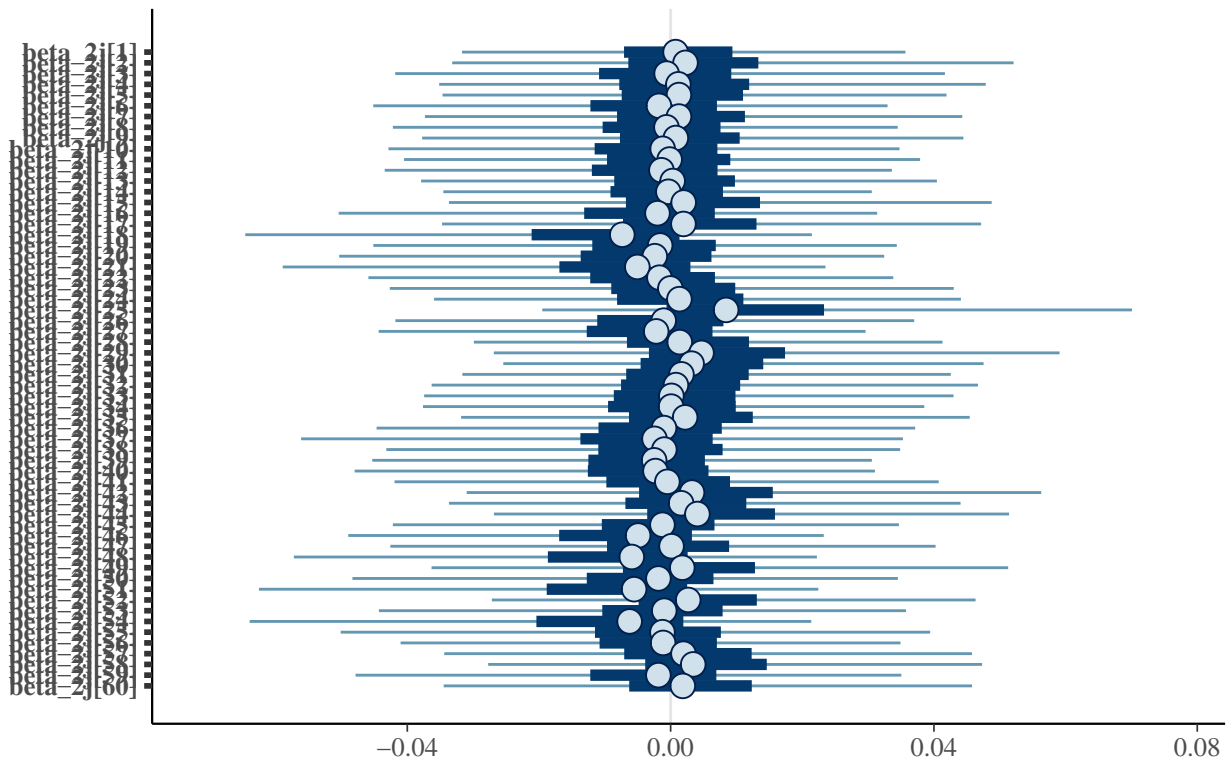
Posterior distributions of coef for urban  
with 95% intervals



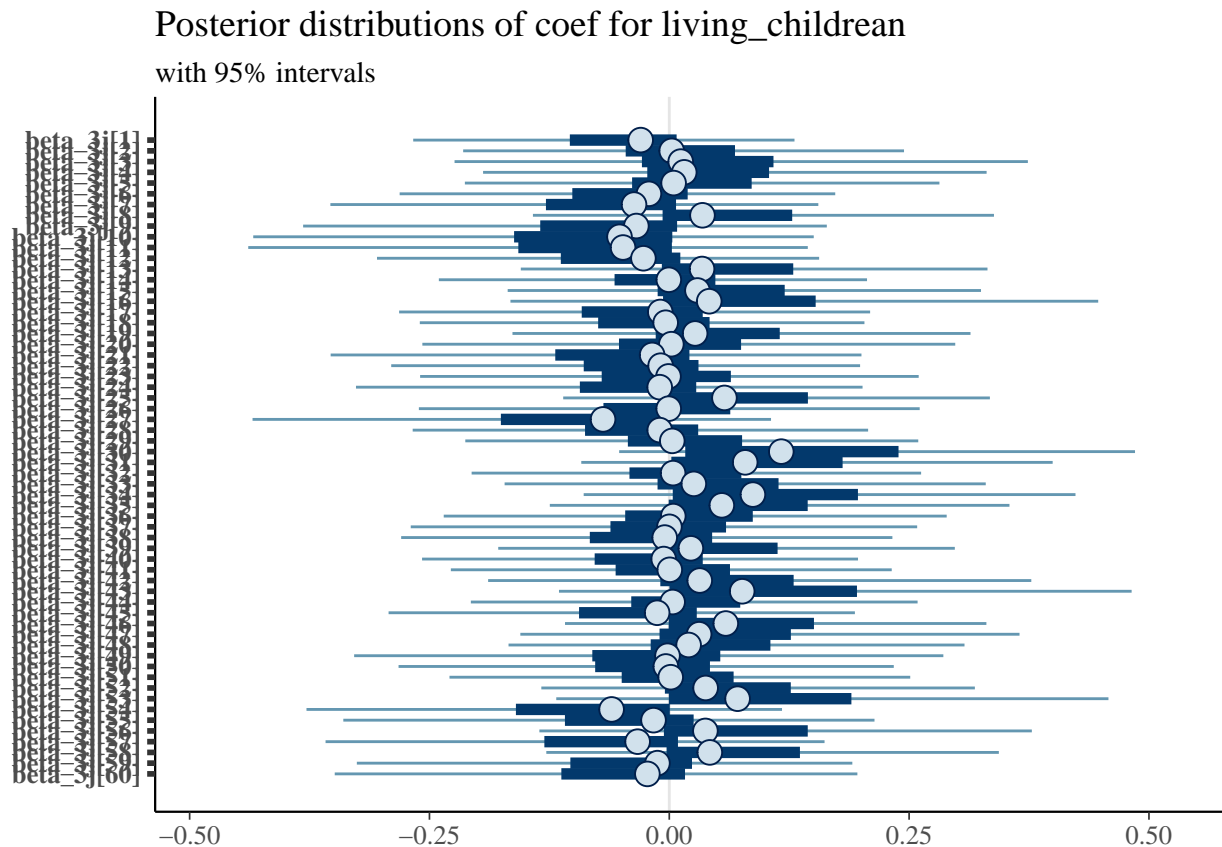
```
mcmc_intervals(fit_posterior, pars=beta_2s, prob_out=0.95) +
  ggplot2::labs(
    title = "Posterior distributions of coef for age_mean",
    subtitle = "with 95% intervals"
  )
```



Posterior distributions of coef for age\_mean  
with 95% intervals



```
mcmc_intervals(fit_posterior, pars=beta_3s, prob_out=0.95) +  
  ggplot2::labs(  
    title = "Posterior distributions of coef for living_childrean",  
    subtitle = "with 95% intervals"  
  )
```



**ANSWER (4. Varying-Coefficients Model (a)):**

Based on the 95% central posterior intervals for the coefficients within each district, coefficients for 'urban' shows greater variation amongst the districts compared to those for 'age\_mean' and 'living.children'. It suggests that the coefficients for 'urban' vary by district while those for 'age\_mean' and 'living.children' do not.

- (b) Use all of the information you've gleaned thus far to build a final Bayesian logistic regression classifier on the train set. Then, use your model to make predictions on the test set. Report your model's classification percentage.

```
# stan code for problem 4
stan_code4b <- c("
data {
  // Number of observations
  int N;
  // Number of districts
  int J;
  // List of features, one for each observation
  int district[N];
  int<lower=0, upper=1> urban[N];
  int living_children[N];
  real age_mean[N];
  // Binary response (integer array)
  int<lower=0, upper=1> Y[N];

  int N_test;
  // Number of districts
  int J_test;
```

```

// List of features, one for each observation
int district_test[N_test];
int<lower=0, upper=1> urban_test[N_test];
int living_children_test[N_test];
real age_mean_test[N_test];
}

parameters {
  real mu_0;
  real<lower=0> sigma_0;
  real beta_0j[J]; // bias term vary by J districts
  real beta_1j[J];
  real<lower=0> sigma_1;
  real beta_2;
  real beta_3;
}

model {
  // Prior
  mu_0 ~ normal(0,100);
  sigma_0 ~ exponential(0.1);
  sigma_1 ~ exponential(0.1);
  beta_2 ~ normal(0, 100);
  beta_3 ~ normal(0, 100);

  // J different beta_0j, beta_1j priors
  for (j in 1:J) {
    beta_0j[j] ~ normal(mu_0, sigma_0);
    beta_1j[j] ~ normal(0,sigma_1);
  }

  // Likelihood
  for (n in 1:N) {
    Y[n] ~ bernoulli_logit(beta_0j[district[n]] +
                          beta_1j[district[n]]*urban[n] + beta_2*age_mean[n] +
                          beta_3*living_children[n]);
  }
}

generated quantities {
  int y_rep[N_test]; // Draws from posterior predictive dist

  for (n in 1:N_test) {
    y_rep[n] = bernoulli_rng(inv_logit(beta_0j[district_test[n]] +
                                       beta_1j[district_test[n]]*urban_test[n] +
                                       beta_2*age_mean_test[n] + beta_3*living_children_test[n]));
  }
}

")

```

```

# fit the model
options(mc.cores = parallel::detectCores())

```

```
fit4b <- stan(model_code = stan_code4b,
  data = stan_list3e,
  iter = 2000,
  chains = 4,
  seed = 46,
  refresh = FALSE)
```

```
## In file included from filef9241cb27e96.cpp:8:
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/BH/include/boost/math/
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/BH/include/boost/confi
## /Library/Frameworks/R.framework/Versions/3.4/Resources/library/BH/include/boost/config/compiler/clang.hpp:
## # define BOOST_NO_CXX11_RVALUE_REFERENCES
##      ^
## <command line>:6:9: note: previous definition is here
## #define BOOST_NO_CXX11_RVALUE_REFERENCES 1
##      ^
## In file included from filef9241cb27e96.cpp:8:
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/RcppEigen/include/Eige
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/RcppEigen/include/Eige
## /Library/Frameworks/R.framework/Versions/3.4/Resources/library/RcppEigen/include/Eigen/src/Core/util/Reena
## #pragma clang diagnostic pop
##      ^
## In file included from filef9241cb27e96.cpp:8:
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/RcppEigen/include/Eige
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/RcppEigen/include/Eige
## /Library/Frameworks/R.framework/Versions/3.4/Resources/library/RcppEigen/include/Eigen/src/Core/util/Reena
## #pragma clang diagnostic pop
##      ^
## In file included from filef9241cb27e96.cpp:8:
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
```







```

## BOOST_MOVE_STD_NS_BEG
## ^
## /Library/Frameworks/R.framework/Versions/3.4/Resources/library/BH/include/boost/move/detail/std_ns_begin.h
## #define BOOST_MOVE_STD_NS_BEG _LIBCPP_BEGIN_NAMESPACE_STD
## ^
## /Library/Developer/CommandLineTools/usr/include/c++/v1/__config:390:52: note: expanded from macro '_LIBCPP_BEGIN_NAMESPACE_STD'
## #define _LIBCPP_BEGIN_NAMESPACE_STD namespace std {inline namespace _LIBCPP_NAMESPACE {
## ^
## In file included from filef9241cb27e96.cpp:8:
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/stan/math/rev/core/set_zero_all_adjoints.hpp:1:
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/stan/math/rev/core/set_zero_all_adjoints.hpp:1:
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/stan/math/rev/core/set_zero_all_adjoints.hpp:1:
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/stan/math/rev/core/set_zero_all_adjoints.hpp:1:
## /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/stan/math/rev/core/set_zero_all_adjoints.hpp:1:
## static void set_zero_all_adjoints() {
## ^
## In file included from filef9241cb27e96.cpp:8:
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/stan/math/rev/core/set_zero_all_adjoints.hpp:1:
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/stan/math/rev/core/set_zero_all_adjoints.hpp:1:
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/stan/math/rev/core/set_zero_all_adjoints.hpp:1:
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/stan/math/rev/core/set_zero_all_adjoints.hpp:1:
## /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/stan/math/rev/core/set_zero_all_adjoints.hpp:1:
## static void set_zero_all_adjoints_nested() {
## ^
## In file included from filef9241cb27e96.cpp:8:
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/stan/math/rev/core/set_zero_all_adjoints.hpp:1:
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/stan/math/rev/core/set_zero_all_adjoints.hpp:1:
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/stan/math/rev/core/set_zero_all_adjoints.hpp:1:
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/stan/math/rev/core/set_zero_all_adjoints.hpp:1:
## /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/stan/math/prim/mat/fun/fft_next_good_size.hpp:1:
## size_t fft_next_good_size(size_t N) {
## ^
## 19 warnings generated.

```

```

y_rep4b_mean <- colMeans(as.matrix(fit4b, pars="y_rep"))
y_pred4b <- y_rep4b_mean >= 0.5
y_true <- test$contraceptive_use
accuracy4b <- sum(y_pred4b == y_true)/length(y_true)
accuracy4b # CHECK >= 0.5 (0.6897622) OR > 0.5 (0.688728)

```

```
## [1] 0.6463289
```

#### ANSWER (4. Varying-Coefficients Model (b)):

Since the 95% central posterior intervals for the intercept and coefficients vary more by district than the coefficients for ‘age\_mean’ and ‘living.children’, we would vary coefficients for intercept and ‘urban’ by district. This model gives an accuracy score of 0.6463289 on test data.

### 5. “Bayesball” (AC 209b students only)

In Major League Baseball (MLB), each team plays 162 games in the regular season. The data in Bayesball.txt contains information about every regular season game in 2017. The data can be read in by the command

```
games2017 = read.table("Bayesball.txt", sep=',')
```

The relevant columns of the data are



```
data2017 = games2017[,c(7,4,11,10)]
names(data_2017) = c("Home", "Away", "Home_Score", "Away_Score")
```

Because we are going to focus on wins versus losses, you will need to convert these data into binary outcomes.

Under the Bradley-Terry model, each team  $i$  is assumed to have some underlying talent parameter  $\pi_i$ . The model states that the probability that team  $i$  defeats opponent  $j$  in any game is:

$$\Pr(\text{team } i \text{ defeats team } j) = \frac{\pi_i}{\pi_i + \pi_j}.$$

where  $i, j \in (1, 2, \dots, 30)$ , since there are 30 teams in MLB. The parameter  $\pi_i$  is team  $i$ 's "strength" parameter, and is required to be positive.

If we let  $V_{ij}$  be the number of times in a season that team  $i$  defeats team  $j$ , and  $n_{ij}$  to be the number of games between them, an entire season of MLB can be described with the following density:

$$p(V | \pi) = \prod_{i=1}^{N-1} \prod_{j=i+1}^N \binom{n_{ij}}{V_{ij}} \left( \frac{\pi_i}{\pi_i + \pi_j} \right)^{V_{ij}} \left( \frac{\pi_j}{\pi_i + \pi_j} \right)^{V_{ji}}$$

Team  $i$ 's victories against team  $j$  follows a binomial distribution governed by the Bradley-Terry probability with the given strength parameters.

Rather than work with the  $\pi_i$ , we will transform the model by letting  $\lambda_i = \log \pi_i$  for all  $i$ . Thus the probability  $i$  defeats  $j$  is

$$\Pr(\text{team } i \text{ defeats team } j) = \frac{\exp(\lambda_i)}{\exp(\lambda_i) + \exp(\lambda_j)}.$$

The advantage of this parameterization is that the  $\lambda_i$  are unconstrained real-valued parameters.

We now carry out a Bayesian analysis of the Bradley-Terry model. We will assume a hierarchical normal prior distribution on the  $\lambda_i$ , that is

$$\lambda_i \sim N(0, \sigma).$$

for all  $i = 1, \dots, N$ . We will also assume that the standard deviation  $\sigma$  has a uniform prior distribution,

$$\sigma \sim \text{Uniform}(0, 50)$$

with a maximum value of 50.

Thus the full model is:

Prior distribution:

$\lambda_i | \sigma \sim N(0, \sigma)$ , with  $\sigma \sim \text{Uniform}(0, 50)$

Model for data:  $V_{ij} | \lambda_i, \lambda_j, n_{ij} \sim \text{Binomial}\left(n_{ij}, \frac{\exp(\lambda_i)}{\exp(\lambda_i) + \exp(\lambda_j)}\right)$

```
games2017 = read.table("data/Bayesball.txt", sep=',')
data2017 = games2017[,c(7,4,11,10)]
names(data2017) = c("Home", "Away", "Home_Score", "Away_Score")
head(data2017)
```

```
##   Home Away Home_Score Away_Score
## 1  ARI  SFN          6          5
## 2  SLN  CHN          4          3
## 3  TBA  NYA          7          3
## 4  CIN  PHI          3          4
## 5  LAN  SDN         14          3
## 6  MIL  COL          5          7
```

(a) Why does this prior distribution on the  $\lambda_i$  and  $\sigma$  make sense? Briefly explain.

**ANSWER (5. Bayesball (a)):**

The prior on  $\lambda_i \sim N(0, \sigma)$  makes sense because it reflects an expected wining probability between 2 teams to be

$$E[\Pr(\text{team } i \text{ defeats team } j)] = \frac{\exp(0)}{\exp(0) + \exp(0)} = \frac{1}{2}.$$

The prior on  $\sigma \sim \text{Uniform}(0, 50)$  makes sense because it gives the standard deviation of the log scale team strength a reasonable limited range with equal probability spreading over.

(b) Implement the model in Stan.

```
n_teams <- length(unique(data2017$Home))
N_ij <- matrix(rep(0, n_teams*n_teams), ncol = n_teams)
V_ij <- matrix(rep(0, n_teams*n_teams), ncol = n_teams)
for (i in 1:length(levels(data2017$Home))) {
  home <- levels(data2017$Home)[i]
  for (j in 1:length(levels(data2017$Away))) {
    away <- levels(data2017$Away)[j]
    if (away != home) {
      df1_ij <- data2017[data2017$Home == home & data2017$Away == away, ]
      df2_ij <- data2017[data2017$Home == away & data2017$Away == home, ]
      N_ij[i, j] <- nrow(df1_ij) + nrow(df2_ij)
      v1 <- nrow(df1_ij[df1_ij$Home_Score > df1_ij$Away_Score, ])
      v2 <- nrow(df2_ij[df2_ij$Home_Score < df2_ij$Away_Score, ])
      V_ij[i, j] <- v1 + v2
    }
    else {
      N_ij[i, j] <- 0
      V_ij[i, j] <- 0
    }
  }
}
```

```
# create list
stan_list_mlb <- list()
stan_list_mlb$n_team <- n_teams
stan_list_mlb$N <- N_ij
stan_list_mlb$V <- V_ij

# stan code
stan_code_mlb <- c("data {
  int n_team; // Number of teams
  int N[n_team, n_team];
  int V[n_team, n_team]; // Bernoulli Response: {0, 1} array
}

parameters {
  real<lower=0> sigma;
  real lambda[n_team];
}

model {
  // Prior
  sigma ~ uniform(0, 50);
  for (t in 1:n_team) {
```

```

    lambda[t] ~ normal(0, sigma);
  }

  // Likelihood
  for (i in 1:n_team) {
    for (j in 1:n_team) {
      V[i, j] ~ binomial(N[i, j], exp(lambda[i])/(exp(lambda[i]) + exp(lambda[j])));
    }
  }
}

generated quantities {
  int V_rep[n_team, n_team];      // Draws from posterior predictive dist
  for (i in 1:n_team) {
    for (j in 1:n_team) {
      V_rep[i, j] = binomial_rng(N[i, j], exp(lambda[i])/(exp(lambda[i]) + exp(lambda[j])));
    }
  }
}

")

```

```

options(mc.cores = parallel::detectCores())
mlb_fit <- stan(model_code = stan_code_mlb,
               data = stan_list_mlb,
               iter = 2000,
               chains = 4,
               seed = 46,
               refresh = FALSE)

```

```

## In file included from filef92467d51dc3.cpp:8:
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/BH/include/boost/math/
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/BH/include/boost/confi
## /Library/Frameworks/R.framework/Versions/3.4/Resources/library/BH/include/boost/config/compiler/clang.hpp:
## # define BOOST_NO_CXX11_RVALUE_REFERENCES
##      ^
## <command line>:6:9: note: previous definition is here
## #define BOOST_NO_CXX11_RVALUE_REFERENCES 1
##      ^
## In file included from filef92467d51dc3.cpp:8:
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/sr
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/RcppEigen/include/Eige

```





```

## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/BH/include/boost/lexic
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/BH/include/boost/conta
## /Library/Frameworks/R.framework/Versions/3.4/Resources/library/BH/include/boost/container/detail/std_fwd.h
## BOOST_MOVE_STD_NS_BEG
## ^
## /Library/Frameworks/R.framework/Versions/3.4/Resources/library/BH/include/boost/move/detail/std_ns_begin.h
## #define BOOST_MOVE_STD_NS_BEG _LIBCPP_BEGIN_NAMESPACE_STD
## ^
## /Library/Developer/CommandLineTools/usr/include/c++/v1/__config:390:52: note: expanded from macro '_LIBCPP_E
## #define _LIBCPP_BEGIN_NAMESPACE_STD namespace std {inline namespace _LIBCPP_NAMESPACE {
## ^
## In file included from filef92467d51dc3.cpp:8:
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/sr
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/RcppEigen/include/Eige
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/RcppEigen/include/Eige
## /Library/Frameworks/R.framework/Versions/3.4/Resources/library/RcppEigen/include/Eigen/src/Core/util/Reena
## #pragma clang diagnostic pop
## ^
## In file included from filef92467d51dc3.cpp:8:
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/sr
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/RcppEigen/include/Eige
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/RcppEigen/include/Eige
## /Library/Frameworks/R.framework/Versions/3.4/Resources/library/RcppEigen/include/Eigen/src/Core/util/Reena
## #pragma clang diagnostic pop
## ^
## In file included from filef92467d51dc3.cpp:8:
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/sr
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/RcppEigen/include/Eige
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/RcppEigen/include/Eige
## /Library/Frameworks/R.framework/Versions/3.4/Resources/library/RcppEigen/include/Eigen/src/Core/util/Reena
## #pragma clang diagnostic pop
## ^
## In file included from filef92467d51dc3.cpp:8:
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/sr
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/RcppEigen/include/Eige
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/RcppEigen/include/Eige
## /Library/Frameworks/R.framework/Versions/3.4/Resources/library/RcppEigen/include/Eigen/src/Core/util/Reena
## #pragma clang diagnostic pop

```

```

##                                     ^
## In file included from filef92467d51dc3.cpp:8:
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/RcppEigen/include/Eige
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/RcppEigen/include/Eige
## /Library/Frameworks/R.framework/Versions/3.4/Resources/library/RcppEigen/include/Eigen/src/Core/util/Reena
##     #pragma clang diagnostic pop
##                                     ^
## In file included from filef92467d51dc3.cpp:460:
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/rstan/include/rstan/rs
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/rstan/include/rstan/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/sr
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/sr
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/sr
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/BH/include/boost/circu
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/BH/include/boost/circu
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/BH/include/boost/move/
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/BH/include/boost/move/
## /Library/Frameworks/R.framework/Versions/3.4/Resources/library/BH/include/boost/move/detail/iterator_trait
## BOOST_MOVE_STD_NS_BEG
## ^
## /Library/Frameworks/R.framework/Versions/3.4/Resources/library/BH/include/boost/move/detail/std_ns_begin.h
##     #define BOOST_MOVE_STD_NS_BEG _LIBCPP_BEGIN_NAMESPACE_STD
##                                     ^
## /Library/Developer/CommandLineTools/usr/include/c++/v1/__config:390:52: note: expanded from macro '_LIBCPP_F
## #define _LIBCPP_BEGIN_NAMESPACE_STD namespace std {inline namespace _LIBCPP_NAMESPACE {
##                                     ^
## In file included from filef92467d51dc3.cpp:8:
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/sr
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/stan/math/rev/core/set_
##     static void set_zero_all_adjoints() {
##                                     ^
## In file included from filef92467d51dc3.cpp:8:
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/sr
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/stan/math/rev/core/set_
##     static void set_zero_all_adjoints_nested() {
##                                     ^
## In file included from filef92467d51dc3.cpp:8:
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/sr
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## In file included from /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/st
## /Library/Frameworks/R.framework/Versions/3.4/Resources/library/StanHeaders/include/stan/math/prim/mat/fun/
##     size_t fft_next_good_size(size_t N) {

```

```
##          ^
## 19 warnings generated.
```

(c) Report the posterior means for each team's exponentiated strength parameters (that is,  $\exp(\lambda_i)$ ).

```
mlb_pi <- exp(as.matrix(mlb_fit, pars = "lambda"))
mean_team_strength <- colMeans(mlb_pi)
names(mean_team_strength) <- levels(data2017$Home)
print(sort(mean_team_strength))
```

```
##      PHI      SFN      DET      NYN      CIN      CHA      ATL
## 0.7183123 0.7210040 0.7389480 0.7759817 0.7760767 0.7885896 0.8023738
##      SDN      PIT      MIA      OAK      BAL      TOR      TEX
## 0.8275531 0.8753876 0.8835597 0.9144326 0.9346382 0.9476207 0.9664571
##      SEA      KCA      ANA      TBA      SLN      MIL      MIN
## 0.9706398 1.0059329 1.0168632 1.0256484 1.0262627 1.0872983 1.1057144
##      COL      CHN      NYA      ARI      WAS      BOS      HOU
## 1.1285521 1.2104051 1.2572780 1.2704498 1.2939080 1.3090778 1.5213175
##      CLE      LAN
## 1.5340730 1.5662956
```

(d) Using the posterior predictive distribution for the strengths of the Dodgers and Astros, simulate 1000 recreations of the 2017 World Series. That is, simulate 1000 separate series between the teams, where the series ends when either team gets to 4 wins. Based on your simulation, what was the probability that the Astros would have won the World Series last year?

```
dodgers_strength <- mean_team_strength[levels(data2017$Home) == 'LAN']
astros_strength <- mean_team_strength[levels(data2017$Home) == 'HOU']
total_dodgers_win <- 0
total_astros_win <- 0
for (i in 1:1000) {
  dodgers_win <- 0
  astros_win <- 0
  while (dodgers_win < 4 & astros_win < 4) {
    dodgers <- rbinom(n=1, size=1, dodgers_strength/(dodgers_strength+astros_strength))
    if (dodgers == 1) {
      dodgers_win <- dodgers_win + 1
    }
    else {
      astros_win <- astros_win + 1
    }
  }
  if (dodgers_win > astros_win) {
    total_dodgers_win <- total_dodgers_win + 1
  }
  else {
    total_astros_win <- total_astros_win + 1
  }
}
```

```
print(paste0('Based on 1000 simulations, the probability that the Astros would have won the 2017 World Series is ', total_astros_win/1000, '%'))
```

```
## [1] "Based on 1000 simulations, the probability that the Astros would have won the 2017 World Series is 48.3 %"
```