

# Assignment 4

## Introduction to Bayesian Data Analysis 2025

Ana Sofia Acevedo (03805441)

### Preamble

- **Points:** Assignment 4 comprises of 5 tasks, 2 points for Tasks 1.1-2.2 each and 4 points for Task 3 (12 in total). Full points are obtained for complete and correct answers—a proportion of the points are obtained for a proper approach or if only part of the task is solved.
- **Submission:** Hand in the assignment as a PDF `Markdown` report. The report should show the results, the code that produced the results, and additional text or comment. The report should appear clean and be uploaded on Moodle until Wednesday, July 7, 9:45 am.
- **Collaboration:** Reports can be handed in as team work (max. 2 people). When working in teams, declare this on page 2. However, each collaborator needs to hand in a report via Moodle, stating their name, student number (p. 1), and their machine specification (p. 2).
- **Permitted and Prohibited:** You may use materials from this class (e.g., slides, code on GitHub) as well as online forums such as [Stack Overflow](#) to write your code. However, you are not allowed to post questions from the assignment online or prompt them (including paraphrases) to LLMs/chatbots. All use of LLMs/chatbots is generally not allowed. Solutions may not be shared with other students from the class (except 1 potential collaborator).

## ulam and Quarto

Running MCMC with `ulam()` takes some time and produces many messages. This can be annoying when you repeatedly render the **Quarto** document and the goal is to have a clean report. Here are two tips to avoid an ugly document and refitting the model every time you re-render the document.

- 1) Write the `ulam()` model in a separate code chunk with the `echo`, `eval`, `output` settings as below to avoid printing the MCMC progress messages of `ulam()` in the PDF document.
- 2) Use the `file` argument in `ulam()` to fit the model only once and store the fitting results for reuse. In the code below, `ulam()` first looks into the folder `assignment_4` for a file `m1.rds`—the model fit object `m1`. If the file exists, it loads this file without fitting the model again. If the file does not exist yet, `ulam()` fits the model and stores the fitted object in the respective folder, so it can be reused later.

```
#| echo: true
#| eval: true
#| output: false
```

```
m1 <- ulam(
  model ,
  data= ... ,
  chains = ... ,
  cores = ... ,
  file = here("assignment_4", "m1")
)
```

## Authorship Information

### 1. Declaration of Collaboration

- ☐ Yes (Collaborator name)
- ☒ No

### 2. Declaration of Authorship

- ☒ I certify that this assignment represents my own work. I have not used any unauthorized or unacknowledged aids as stated in the preamble, including free or commercial systems or services offered on the internet or text generating systems embedded into software. I did not copy code from someone else nor did I share my code with someone else.

### 3. System Information

- ☒ I confirm that I generated the submitted PDF report myself using R version 4.5.0 (2025-04-11 ucrt) and Quarto/RMarkdown.

Machine stamp: x86\_64-w64-mingw32/x64

Timestamp: 2025-07-15 10:05:15 CEST

```
library(rethinking)
```

```
Loading required package: cmdstanr
```

```
This is cmdstanr version 0.9.0
```

- CmdStanR documentation and vignettes: [mc-stan.org/cmdstanr](https://mc-stan.org/cmdstanr)
- CmdStan path: C:/Users/57314/.cmdstan/cmdstan-2.36.0
- CmdStan version: 2.36.0

```
Loading required package: posterior
```

```
This is posterior version 1.6.1
```

```
Attaching package: 'posterior'
```

```
The following objects are masked from 'package:stats':
```

```
mad, sd, var
```

```
The following objects are masked from 'package:base':
```

```
%in%, match
```

```
Loading required package: parallel
```

```
rethinking (Version 2.42)
```

```
Attaching package: 'rethinking'
```

```
The following object is masked from 'package:stats':
```

```
rstudent
```

```
library(tidyverse)
```

```
-- Attaching core tidyverse packages ----- tidyverse 2.0.0 --  
v dplyr      1.1.4      v readr      2.1.5  
v forcats    1.0.0      v stringr    1.5.1  
v ggplot2     3.5.2      v tibble     3.2.1  
v lubridate  1.9.4      v tidyr      1.3.1  
v purrr       1.0.4
```

```
-- Conflicts ----- tidyverse_conflicts() --  
x dplyr::filter() masks stats::filter()  
x dplyr::lag()     masks stats::lag()  
x purrr::map()     masks rethinking::map()  
i Use the conflicted package (<http://conflicted.r-lib.org/>) to force all conflicts to become
```

```
library(here)
```

here() starts at C:/Users/57314/BayesIntro25\_forked

```
library(ggplot2)
```

## Task Set 1

For Tasks 1.1 and 1.2, consider the data generating processes for two groups of people below. Each of the  $i=1, \dots, N$  ( $N=500$ ) persons belongs to either Group 1 or Group 2 and either succeeded in a task ( $\text{outcome}=1$ ) or not ( $\text{outcome}=0$ ).

```
N <- 500 # number of people  
group <- rbinom(N, 1, prob=.5) # determine group membership  
outcome <- rbinom(N, 1, prob=ifelse(group==0, .2, .5)) # determine outcome  
sim <- data.frame(person=1:N, group, outcome)
```

### Task 1.1

Model the outcome data with a Bernoulli likelihood using an *indicator* variable for group membership. Plot the posterior distribution of the *probability of success* for each group separately.

## Answer

To use a Bernoulli likelihood, we work with the data for each individual.

```
m1_indicator <- ulam(  
  alist(  
    # likelihood  
    outcome ~ dbern(p) ,  
    logit(p) <- a + b*group ,  
  
    # priors  
    a ~ dnorm(0,1.5) ,  
    b ~ dnorm(0,1.5)  
  
  ) ,  
  data=sim,  
  chains = 4,  
  cores = 4,  
  file = here("my_scripts/my_models", "assignment4_m1_indicator")  
)
```

```
# Results analysis  
precis(m1_indicator, depth = 2)
```

	mean	sd	5.5%	94.5%	rhat	ess_bulk
a	-1.257723	0.1409938	-1.4850560	-1.040856	1.002989	691.0131
b	1.220480	0.1847214	0.9364958	1.523236	1.004293	712.3438

```
#alpha (mean probability for group 1)  
inv_logit(coef(m1_indicator)[1])
```

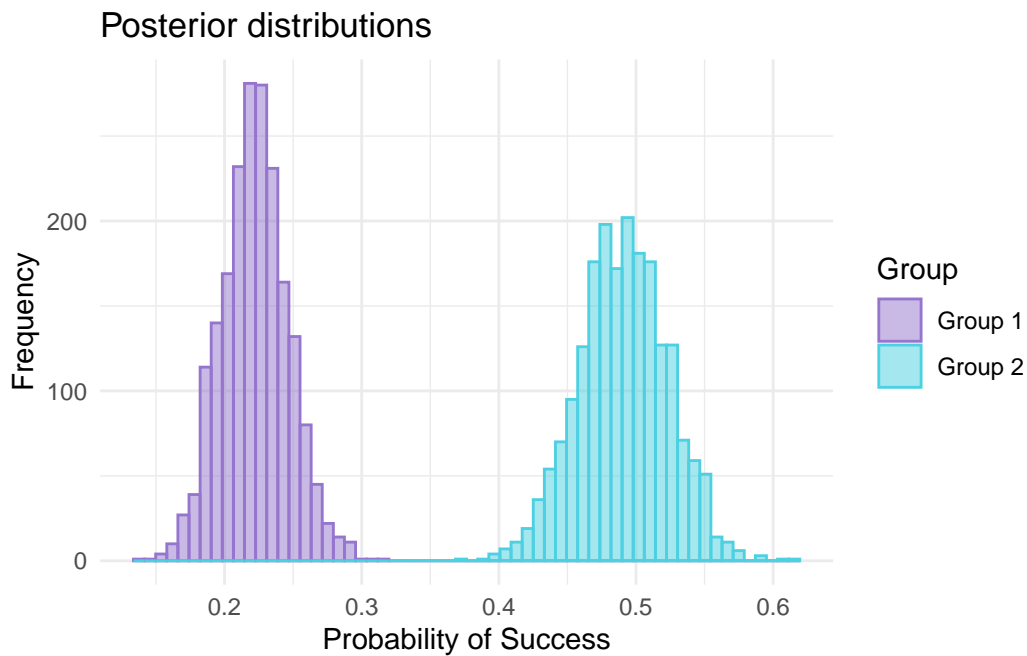
```
      a  
0.2213661
```

```
#alpha+beta (mean probability for group 2)  
inv_logit(sum(coef(m1_indicator)))
```

```
[1] 0.4906902
```

```
samp_m1_indicator <- extract.samples(m1_indicator)
m1_indicator.post <- data.frame(
  post_g1 = inv_logit(samp_m1_indicator$a),
  post_g2 = inv_logit(samp_m1_indicator$a + samp_m1_indicator$b)
)

ggplot() +
  geom_histogram(data = m1_indicator.post, aes(x = post_g1, fill = "Group 1"),
    alpha = 0.5, color = "#9575CD", bins = 60) +
  geom_histogram(data = m1_indicator.post, aes(x = post_g2, fill = "Group 2"),
    alpha = 0.5, color = "#4DD0E1", bins = 60) +
  scale_fill_manual(values = c("Group 1" = "#9575CD", "Group 2" = "#4DD0E1")) +
  labs(x = "Probability of Success",
    y = "Frequency",
    fill = "Group", title = "Posterior distributions") +
  theme_minimal()
```



## Task 1.2

Model the outcome data with a Binomial likelihood using an *index* variable for group membership. Plot the posterior distribution of the *differences in the probability of success* between the

two groups—i.e., the 95% percentile interval should exclude 0 if there is a credible difference.

### Answer

To use a Binomial likelihood, we work with the aggregated data, where the number of trials ‘N’ is included.

```
sim_wide <- list(  
  group = 1:2,  
  outcome = as.integer(tapply(sim$outcome, sim$group, sum)),  
  total = as.integer(tapply(sim$outcome, sim$group, length))  
)  
view(sim_wide)
```

```
m1_index <- ulam(  
  alist(  
    # likelihood  
    outcome ~ dbinom(total, p) ,  
    logit(p) <- a[group] ,  
  
    # priors  
    a[group] ~ dnorm(0,1.5)  
  
  ) ,  
  data=sim_wide,  
  chains = 4,  
  cores = 4,  
  file = here("my_scripts/my_models", "assignment4_m1_index")  
)
```

```
# Result analysis  
precis(m1_index, depth = 2)
```

	mean	sd	5.5%	94.5%	rhat	ess_bulk
a[1]	-1.25402387	0.1528763	-1.5047503	-1.0144372	1.001202	1171.561
a[2]	-0.02536868	0.1292964	-0.2358079	0.1797199	1.000950	1180.641

```
inv_logit(coef(m1_index))
```

	a[1]	a[2]
	0.2220044	0.4936582



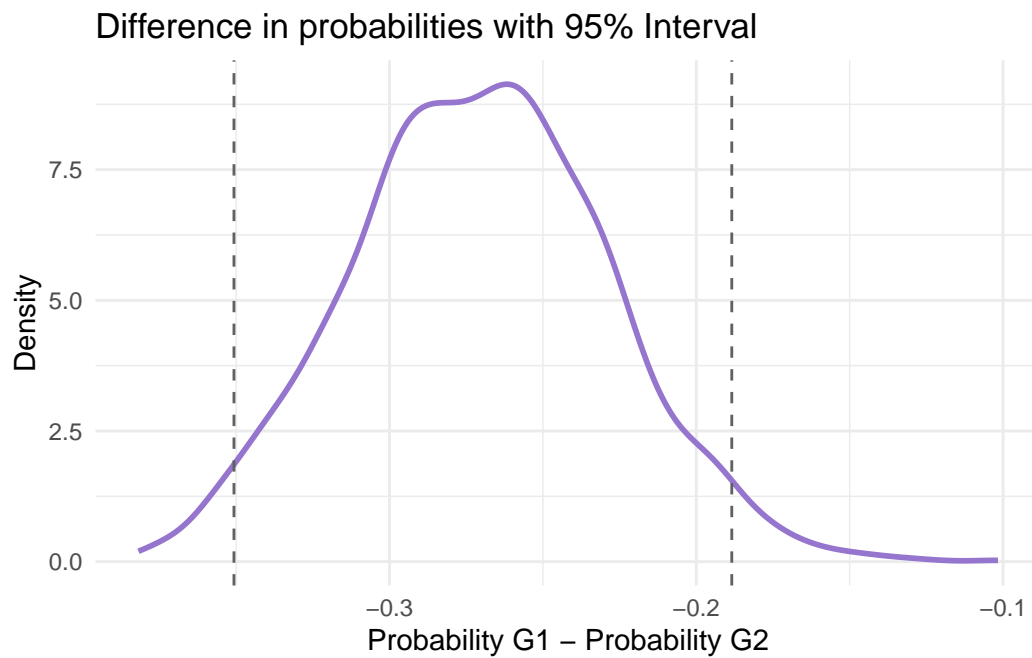
```

m1_index.post <- extract.samples(m1_index)
m1_index.post <- tibble(p1 = plogis(m1_index.post$a[,1]) ,
                        p2 = plogis(m1_index.post$a[,2]) ,
                        diff = p1 - p2)

ci <- quantile(m1_index.post$diff, probs=c(0.025, 0.975))

ggplot(m1_index.post, aes(diff)) +
  geom_density(linewidth = 1, color = "#9575CD") +
  geom_vline(xintercept = ci, linetype = "dashed", color = "#616161") +
  labs(x = "Probability G1 - Probability G2",
       y = "Density",
       title = "Difference in probabilities with 95% Interval") +
  theme_minimal()

```



## Task Set 2

Load the data set `heart.csv` to solve the Task Set 2.

```
heart <- read.csv(here("my_scripts/my_assignments", "heart.csv"))
dat_heart <- data.frame(
  CHD = heart$TenYearCHD,
  # man is 1, woman is 2
  G = ifelse(heart$male==1, 1, 2),
  # 1 the person has diabetes, 2 they dont have
  D = ifelse(heart$diabetes==1, 1, 2),
  A = heart$age
)
```

### Task 2.1

Run a Bayesian logistic regression model to estimate the risk of men and women with and without diabetes to develop a coronary heart disease (TenYearCHD). Does diabetes increase the risk for men or women?

#### Answer

The dependent variable is coronary heart disease, and the independent variables are diabetes and gender.

```
m2_diabetes <- ulam(
  alist(
    # likelihood
    CHD ~ dbern(p) ,
    logit(p) <- a[D, G],

    # priors
    matrix[D,G]:a ~ dnorm(0,1)
  ) ,
  data=dat_heart,
  chains = 4,
  cores = 4,
  file = here("my_scripts/my_models", "assignment4_m2_diabetes")
)
```

```
# analyse the results here
#precis(m2_diabetes, depth=3)
df <- data.frame(case = c("D,M", "ND,M", "D,W",
                          "ND,W"),
                 prob = inv_logit(coef(m2_diabetes)))

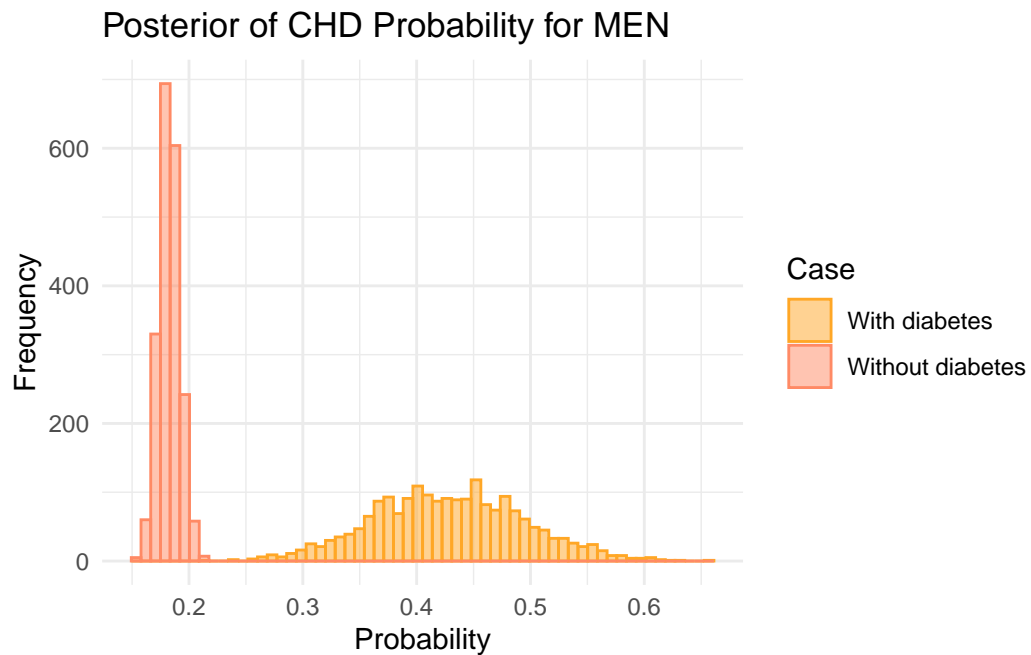
df
```

```
      case      prob
a[1,1]  D,M 0.4265196
a[2,1]  ND,M 0.1825881
a[1,2]  D,W 0.3247633
a[2,2]  ND,W 0.1205424
```

```
#Extracting samples to then calculate the differences and check if there is a
#differential risk
m2_diabetes.post <- extract.samples(m2_diabetes)

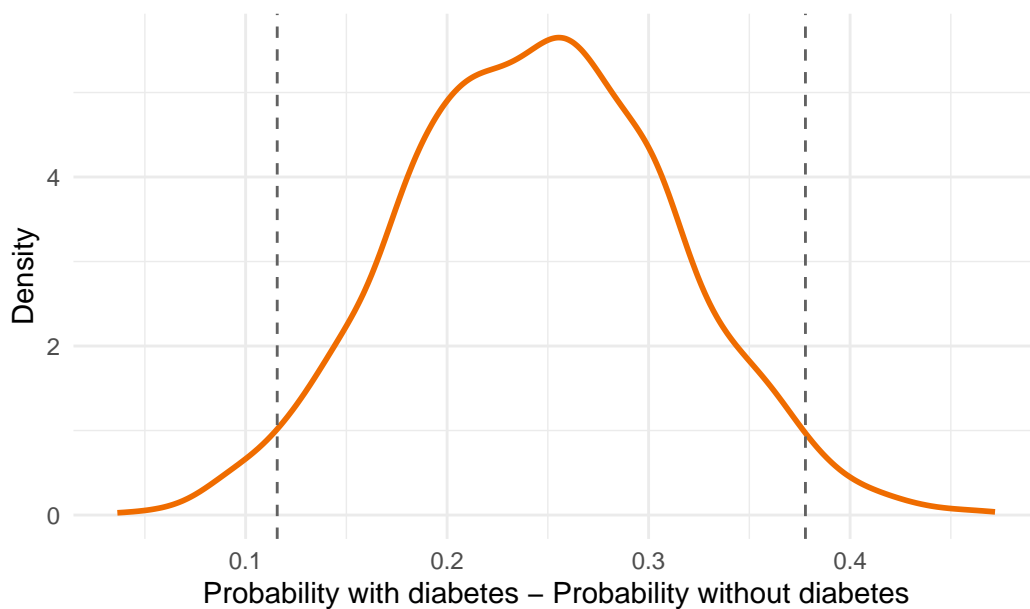
# For MEN
m2_diabetes_men.post <- tibble(p1 = plogis(m2_diabetes.post$a[,1,1]) ,
                              p2 = plogis(m2_diabetes.post$a[,2,1]) ,
                              diff = p1 - p2)
ci_men <- quantile(m2_diabetes_men.post$diff, probs=c(0.025, 0.975))

ggplot() +
  geom_histogram(data = m2_diabetes_men.post,
                aes(x = p1, fill = "With diabetes"),
                alpha = 0.5, color = "#FFA726", bins = 60) +
  geom_histogram(data = m2_diabetes_men.post,
                aes(x = p2, fill = "Without diabetes"),
                alpha = 0.5, color = "#FF8A65", bins = 60) +
  scale_fill_manual(values = c("With diabetes" = "#FFA726",
                              "Without diabetes" = "#FF8A65")) +
  labs(
    x = "Probability",
    y = "Frequency",
    fill = "Case",
    title = "Posterior of CHD Probability for MEN"
  ) +
  theme_minimal()
```



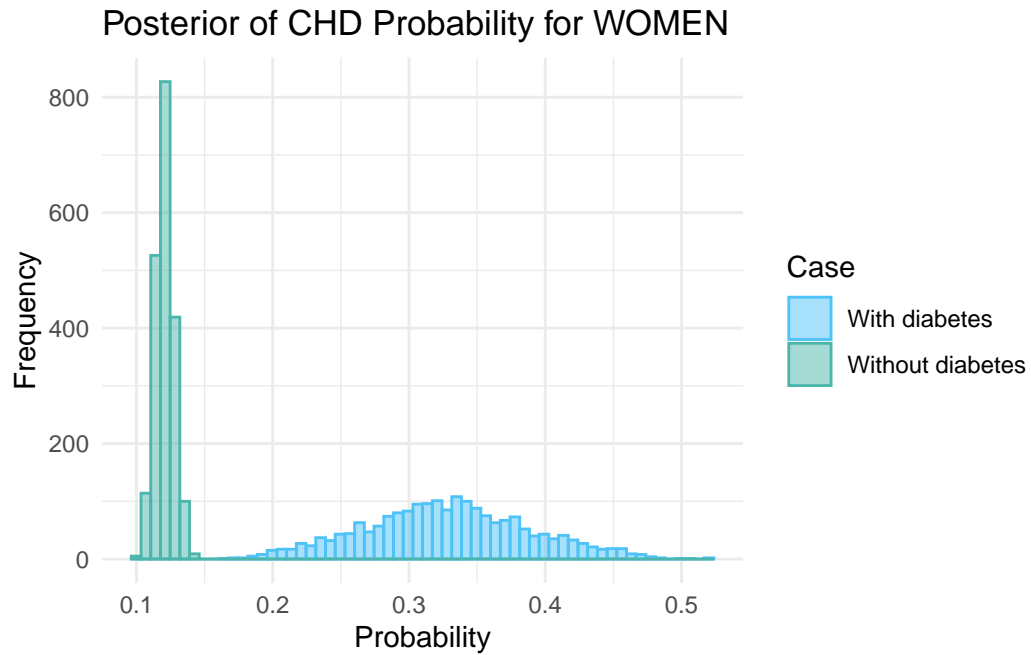
```
ggplot(m2_diabetes_men.post, aes(diff)) +  
  geom_density(linewidth = 1, color = "#EF6C00") +  
  geom_vline(xintercept = ci_men, linetype = "dashed", color = "#616161") +  
  labs(x = "Probability with diabetes - Probability without diabetes",  
       y = "Density",  
       title = "Difference in probabilities for MEN with 95% Interval") +  
  theme_minimal()
```

Difference in probabilities for MEN with 95% Interval

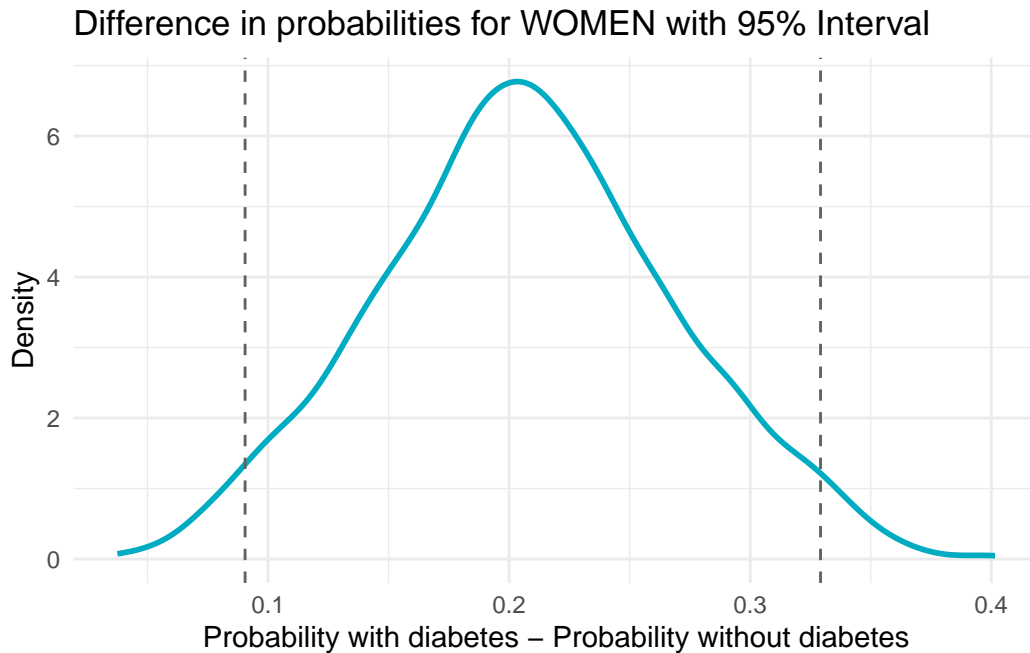


```
# For WOMEN
m2_diabetes_women.post <- tibble(p1 = plogis(m2_diabetes.post$a[,1,2]) ,
                                p2 = plogis(m2_diabetes.post$a[,2,2]) ,
                                diff = p1 - p2)
ci_women <- quantile(m2_diabetes_women.post$diff, probs=c(0.025, 0.975))

ggplot() +
  geom_histogram(data = m2_diabetes_women.post, aes(x = p1,
                                                    fill = "With diabetes"),
                alpha = 0.5, color = "#4FC3F7", bins = 60) +
  geom_histogram(data = m2_diabetes_women.post, aes(x = p2,
                                                    fill = "Without diabetes"),
                alpha = 0.5, color = "#4DB6AC", bins = 60) +
  scale_fill_manual(values = c("With diabetes" = "#4FC3F7",
                              "Without diabetes" = "#4DB6AC")) +
  labs(
    x = "Probability",
    y = "Frequency",
    fill = "Case",
    title = "Posterior of CHD Probability for WOMEN"
  ) +
  theme_minimal()
```



```
ggplot(m2_diabetes_women.post, aes(diff)) +  
  geom_density(linewidth = 1, color = "#00ACC1") +  
  geom_vline(xintercept = ci_women, linetype = "dashed", color = "#616161") +  
  labs(x = "Probability with diabetes - Probability without diabetes",  
       y = "Density",  
       title = "Difference in probabilities for WOMEN with 95% Interval") +  
  theme_minimal()
```



Given that both intervals exclude 0, there is a credible difference in the risk of developing CHD given a preexisting condition of diabetes for both men and women. Specifically, diabetes increases the risk, since the difference of probabilities is positive in both cases. For women, having diabetes increases the risk by around 20% and for men by around 25%.

## Task 2.2

Run a Bayesian logistic regression model to estimate the effect of age on the risk of developing a coronary heart disease (TenYearCHD), separately for women and men. Ensure that the regression intercept represents the risk of women and men with average age. Does age increase the risk for men or women?

## Answer

The dependent variable is coronary heart disease, and the independent variables are age and gender. The model has to be centered around the age variable. I tried using both indicator and index, considering that for both of these the gender variable has to be labeled differently for the package to work.

I also centered the data before, instead of calculating the average and using it inside the `ulam()` function, because I was getting an error when trying to do that. The age is averaged with respect to all the individuals in the sample.

```
dat_heart$A_c <- dat_heart$A - mean(dat_heart$A)
```

```
m2_gender_index <- ulam(  
  alist(  
    # likelihood  
    CHD ~ dbern(p),  
    logit(p) <- a[G]+b_1*(A_c),  
  
    # priors  
    a[G] ~ dnorm(0,1),  
    b_1 ~ dnorm(0,1)  
  
  ) ,  
  data=dat_heart,  
  chains = 4,  
  cores = 4,  
  file = here("my_scripts/my_models", "assignment4_m2_gender_index")  
)  
  
# Changing the labels of the gender variable to be able to implement it as  
# an indicator, gender is 1 if it is a man, 0 if woman, so women is the  
# base group  
dat_heart$G <- ifelse(dat_heart$G==1,1,0)  
  
m2_gender_indicator <- ulam(  
  alist(  
    # likelihood  
    CHD ~ dbern(p) ,  
    logit(p) <- a+b_1*(A_c)+b_2*G,  
  
    # priors  
    a ~ dnorm(0,1),  
    b_1 ~ dnorm(0,1),  
    b_2 ~ dnorm(0,1)  
  
  ) ,  
  data=dat_heart,  
  chains = 4,  
  cores = 4,  
  file = here("my_scripts/my_models", "assignment4_m2_gender_indicator")  
)
```



```
# analyse the results here
# Using the index approach
# traceplot(m2_gender_index)
precis(m2_gender_index, depth = 3)
```

	mean	sd	5.5%	94.5%	rhat	ess_bulk
a[1]	-1.55951771	0.065362261	-1.66766440	-1.45708890	1.003707	1282.290
a[2]	-2.11828453	0.065973393	-2.22438380	-2.01531680	1.001485	1319.654
b_1	0.07609784	0.005177476	0.06783173	0.08442597	1.003961	1599.324

```
#a[1] is the intercept for men, a[2] the intercept for women,
# and b_1 age coefficient
inv_logit(coef(m2_gender_index))
```

	a[1]	a[2]	b_1
	0.1737159	0.1073323	0.5190153

```
# Using the indicator approach
# traceplot(m2_gender_indicator)
precis(m2_gender_indicator, depth = 3)
```

	mean	sd	5.5%	94.5%	rhat	ess_bulk
a	-2.11880618	0.068193374	-2.2291287	-2.01062890	1.002112	793.8951
b_1	0.07610011	0.005359403	0.0676616	0.08482098	1.001175	1598.3923
b_2	0.55118985	0.087655532	0.4104082	0.68765774	1.005311	799.8575

```
# intercept for men
inv_logit(coef(m2_gender_indicator)[1]+coef(m2_gender_indicator)[3])
```

	a
	0.1725565

```
# intercept for women
inv_logit(coef(m2_gender_indicator)[1])
```

	a
	0.1072824

```
# age coefficient  
inv_logit(coef(m2_gender_indicator)[2])
```

```
      b_1  
0.5190159
```

Given the positive value of the coefficient for the variable, age does increase the risk for men and women. For each year above the average age, the risk increases by around 51% for both genders.

## Task 3

Assume that data is randomly generated from a normal distribution according to:

$$y_i \sim N(\mu = 100, \sigma = 10) .$$

Simulate  $N = 500$  data points from this distribution. Write a MCMC (Metropolis) algorithm that fits a normal distribution to the simulated data and recovers the mean and standard deviation. Use 2000 iterations and discard the first 1,000 as warm-up. Plot the posterior distribution and the MCMC chain for both parameters.

### Answer

Taking the analogy we did in class, the “islands” would be the parameter values, the “population sizes” are the (posterior) probabilities at each parameter value, and the “days” are samples taken from the joint posterior of the parameters in the model, in this case I guess the iterations.

Given that we are interested in maximizing the posterior probability, going back to Bayes Theorem, the posterior is proportional to the multiplication of the likelihood and the prior.

Then, the posterior  $\propto$  likelihood \* prior. Or, more detailed,

$$p(\theta|Data) \propto p(Data|\theta) * p(\theta)$$

Moreover, by assuming a normal distribution and taking the parameters as independent, then the joint prior will end up being the multiplication, and the previous expression for the posterior can be rewritten as

$$p(Data|\mu, \sigma) * p(\mu) * p(\sigma)$$

The likelihood function can be calculated according to the probability density function of an assumed distribution, in this case normal. Since the observations are taken to be independent, the likelihood will be the product of the density function for all observed data points.

To do this task, I used the following resources:

[Rampinelli, 2019](#)

[Eisenberg, 2017](#)

```

#simulated data
simu <- rnorm(500, mean=100, sd=10)

# likelihood function supposing the data comes from a normal distribution with
# unknown parameters, since the draws are independent, then the joint
# likelihood is the product of the individual results when evaluation the
# probability density function. Since the numbers are too small, working with
# logarithms is better because the multiplication turns into an addition.
likelihood = function(data, mu, sigma) {
  sum(dnorm(data, mu, sigma, log = TRUE))
}

# Then some priors for both parameters are set for both unknown parameters.
# They mean is assumed to be normally distributed, and the standard deviation
# also but it is shifted a lot to the right to avoid it taking negative values
# as much as possible
# mu ~ N(90, 5)
# sd ~ N(15, 0.5)
mu_curr <- 90#rnorm(1, 90, 5)
sigma_curr <- 12#rnorm(1, 12, 2)

# Running the algorithm
iter <- 2000
param <- data.frame(mu = rep(0,iter), sigma = rep(0,iter),
                    posterior = rep(0,iter))

for (i in c(1:iter)){
  # recording the parameters
  param[i,"mu"] = mu_curr
  param[i,"sigma"] = sigma_curr

  # calculating posterior distribution of the current values of parameters
  likelihood_current <- likelihood(simu, mu_curr, sigma_curr)
  prior_mu_current <- dnorm(mu_curr, mean = 90, sd = 5, log = TRUE)
  prior_sigma_current <- dnorm(sigma_curr, mean = 12, sd = 2, log = TRUE)
  post_current <- likelihood_current + prior_mu_current + prior_sigma_current
  param[i,"posterior"] <- post_current

  # generate new proposal based on the proposed distribution
  mu_new <- rnorm(1, mu_curr, 5)
  sigma_new <- rnorm(1, sigma_curr, 1)
}

```

```

# calculating posterior with the new proposal
likelihood_new <- likelihood(simu, mu_new, sigma_new)
prior_mu_new <- dnorm(mu_new, mean = 90, sd = 5, log = TRUE)
prior_sigma_new <- dnorm(sigma_new, mean = 12, sd = 2, log = TRUE)
post_new <- likelihood_new + prior_mu_new + prior_sigma_new

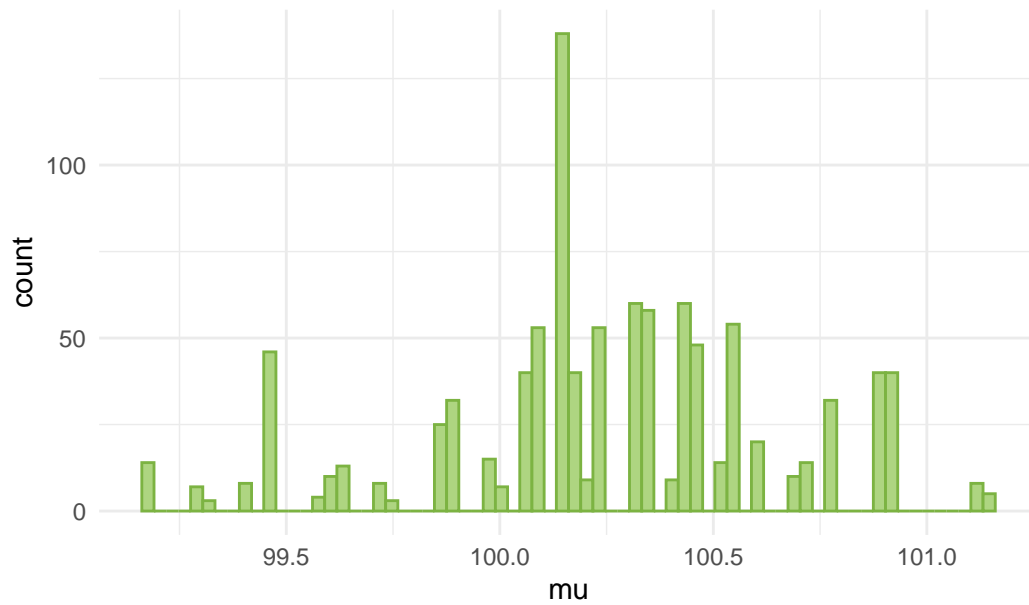
# calculate the ratio to check if the new one is "visited" or not, since it
# is based on the log scale, the ratio is calculated by subtraction. To turn
# it into a comparable ratio with the algorithm it is reversed using "exp"
ratio <- exp(post_new-post_current)
if(ratio>runif(1)){
  mu_curr <- mu_new
  sigma_curr <- sigma_new
}
}

# Ignore the first 1000 as warm-up period for plotting the posterior
post_param <- tail(param, 1000)

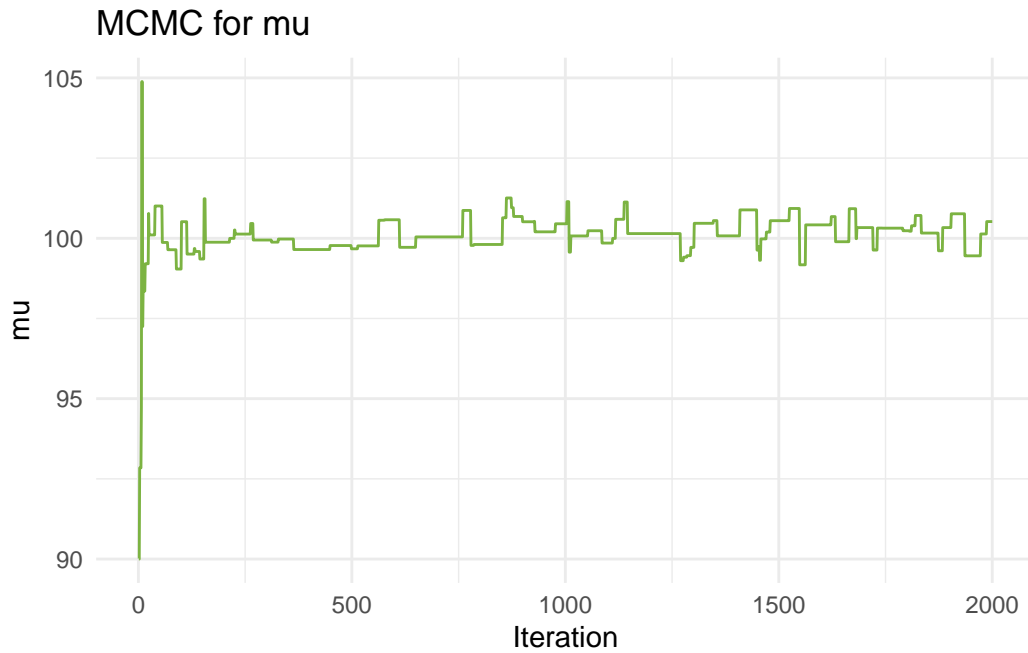
# plot of the posterior probability for mu
ggplot(post_param, aes(x = mu)) +
  geom_histogram(bins = 70, fill = "#AED581", color = "#7CB342") +
  labs(title = "Posterior distribution of mu", x = "mu") +
  theme_minimal()

```

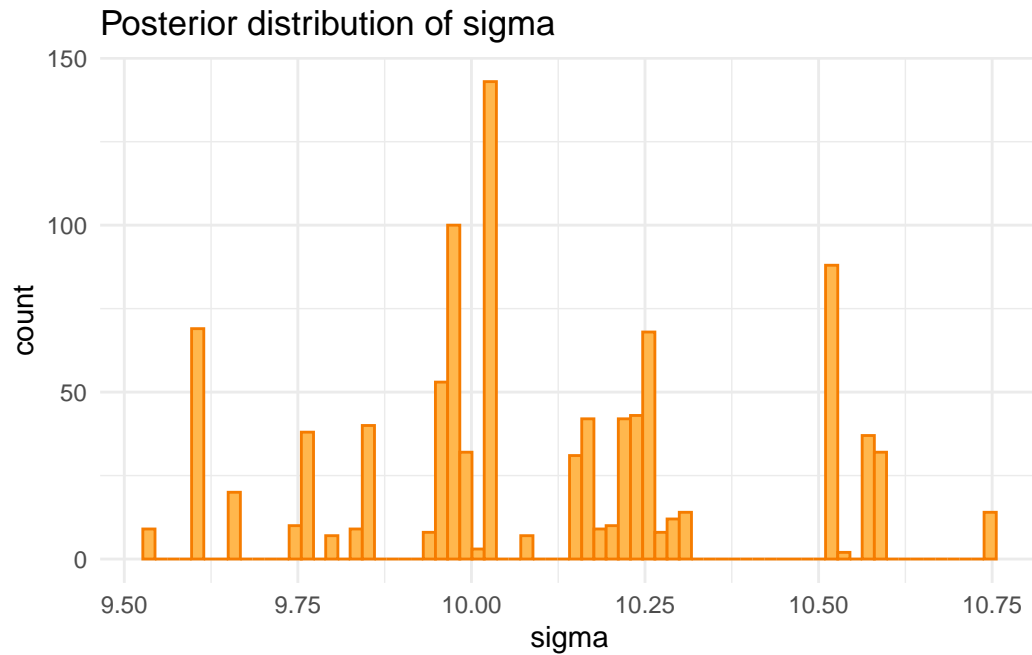
Posterior distribution of mu



```
# Trace plot considering all iterations to show where the chain started and
# where it converged
ggplot(param, aes(x = 1:iter, y = mu)) +
  geom_line(color = "#7CB342") +
  labs(title = "MCMC for mu", x = "Iteration", y = "mu") +
  theme_minimal()
```



```
# plot of the posterior for sigma
ggplot(post_param, aes(x = sigma)) +
  geom_histogram(bins = 70, fill = "#FFB74D", color = "#F57C00") +
  labs(title = "Posterior distribution of sigma", x = "sigma") +
  theme_minimal()
```



```
# Trace plot considering all iterations to show where the chain started and
# where it converged
ggplot(param, aes(x = 1:iter, y = sigma)) +
  geom_line(color = "#F57C00") +
  labs(title = "MCMC for sigma", x = "Iteration", y = "sigma") +
  theme_minimal()
```



