# Lab Exercises 5

16/02/24

## Introduction

Today we will be starting off using Stan, looking at the kid's test score data set (available in resources for the Gelman Hill textbook).

```
library(tidyverse)
library(rstan)
library(tidybayes)
library(here)
```

The data look like this:

```
kidiq <- read_rds(here("kidiq.RDS"))
kidiq
```

```
# A tibble: 434 x 4
   kid_score mom_hs mom_iq mom_age
       <int>  <dbl>  <dbl>   <int>
 1        65      1   121.      27
 2        98      1    89.4     25
 3        85      1   115.      27
 4        83      1    99.4     25
 5       115      1    92.7     27
 6        98      0   108.      18
 7        69      1   139.      20
 8       106      1   125.      23
 9       102      1    81.6     24
10        95      1    95.1     19
# i 424 more rows
```

As well as the kid's test scores, we have a binary variable indicating whether or not the mother completed high school, the mother's IQ and age.
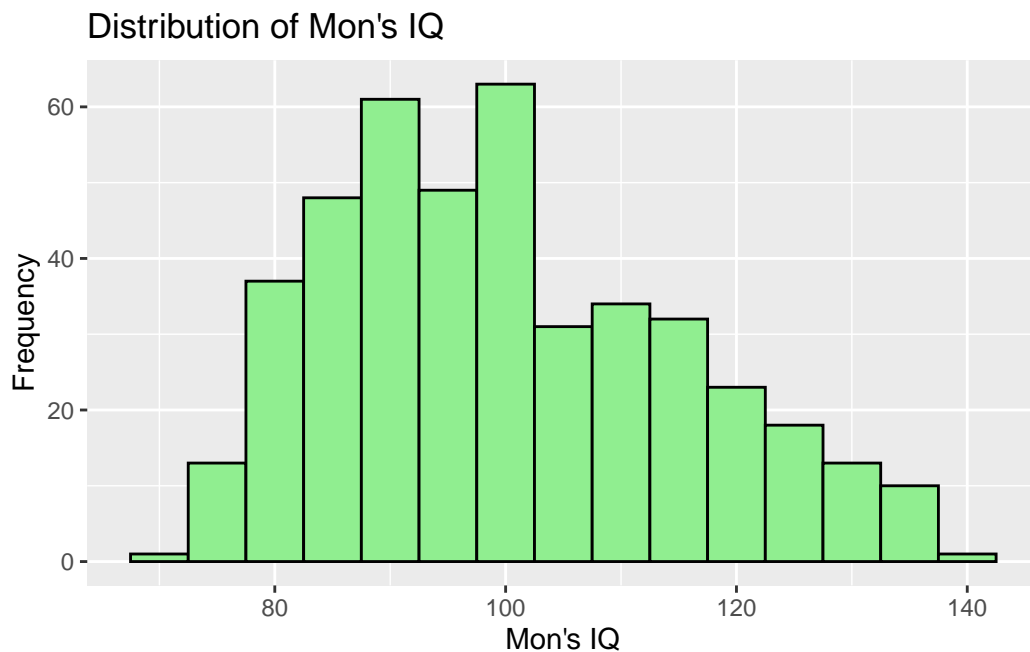
## Descriptives

### Question 1

Use plots or tables to show three interesting observations about the data. Remember:
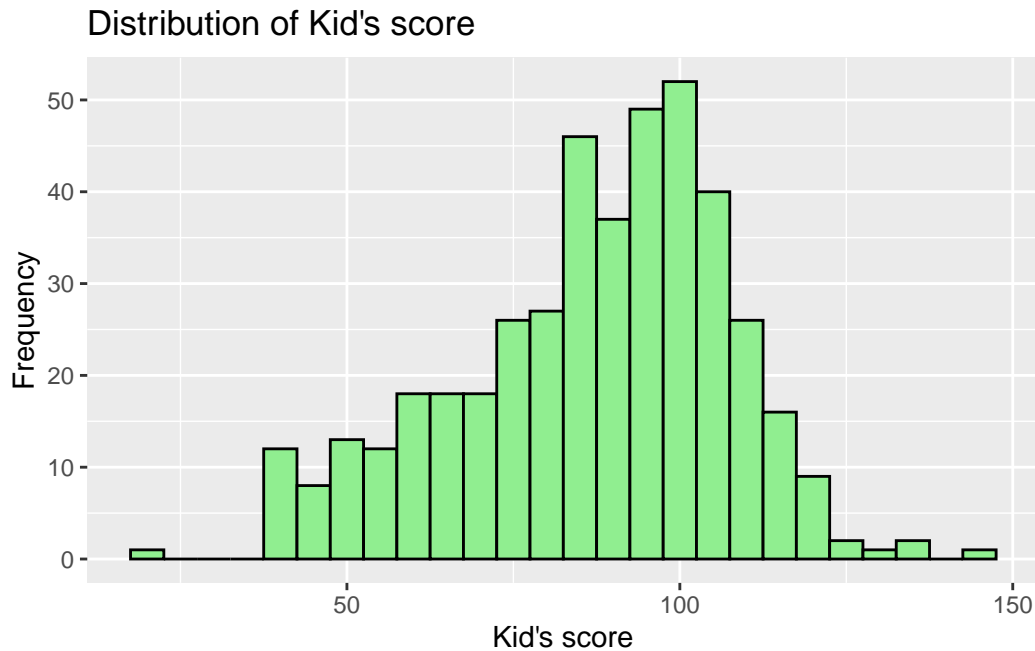
- Explain what your graph/ tables show
- Choose a graph type that's appropriate to the data type

```
ggplot(kidiq, aes(x = mom_iq)) +
  geom_histogram(binwidth = 5, fill = "lightgreen", color = "black") +
  labs(title = "Distribution of Mon's IQ", x = "Mon's IQ", y = "Frequency")
```
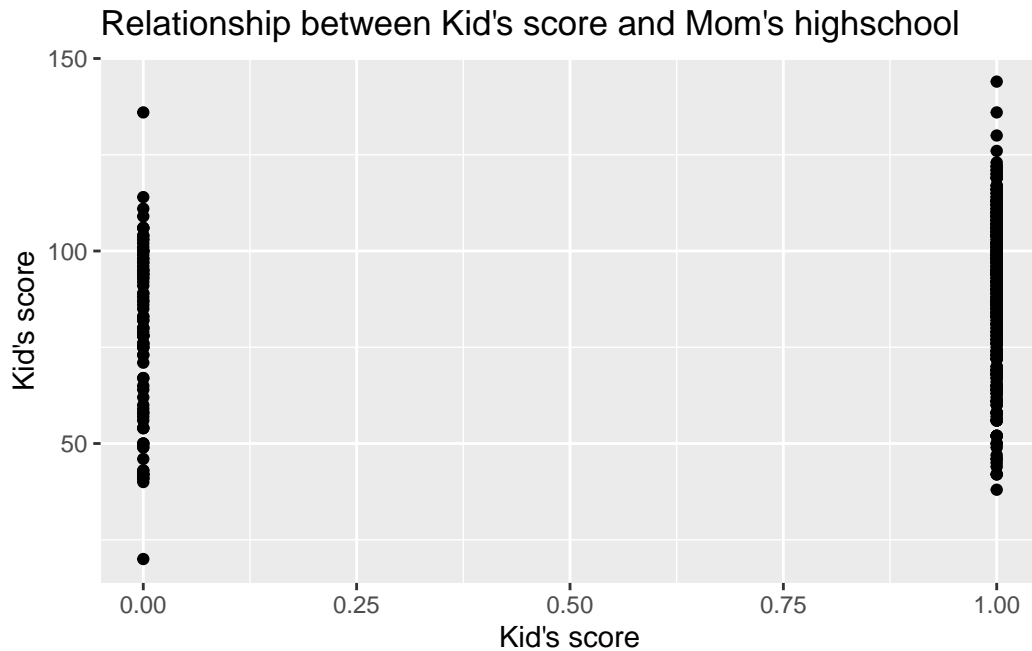


First, the distribution of mom's IQ is explored to get an intuition about the lower, average, and higher scores. The mode is at 100, and the lowest and highest scores are lovated around 70 and 140, respectively.

2

```
ggplot(kidiq, aes(x = kid_score)) +
  geom_histogram(binwidth = 5, fill = "lightgreen", color = "black") +
  labs(title = "Distribution of Kid's score", x = "Kid's score", y = "Frequency")
```

**Distribution of Kid's score**



In the same way of mom's IQ, you can get an intuition of how the kid's score is distributed using the histogram here. The mode is around at 100, and the lowest (ignore outlier) and highest score are around 30 and 130, respectively.

```
ggplot(kidiq, aes(x = mom_hs, y = kid_score)) +
  geom_point() +
  labs(title = "Relationship between Kid's score and Mom's highschool", x = "Kid's score",
```

Relationship between Kid's score and Mom's highschool

It's worth to know how the mom's highschool affects the kid's score as an additional covariate other than the mom's IQ. You will have a slightly higher kid's score when the mom completed highschool, however, it does not have a huge effect on it.

## Estimating mean, no covariates

In class we were trying to estimate the mean and standard deviation of the kid's test scores. The `kids2.stan` file contains a Stan model to do this. If you look at it, you will notice the first `data` chunk lists some inputs that we have to define: the outcome variable y, number of observations N, and the mean and standard deviation of the prior on `mu`. Let's define all these values in a `data` list.

```
y <- kidiq$kid_score
mu0 <- 80
sigma0 <- 10

# named list to input for stan function
data <- list(y = y,
             N = length(y),
             mu0 = mu0,
```

```
                sigma0 = sigma0)
```

Now we can run the model:

```
   fit <- stan(file = here("code/models/kids2.stan"),
               data = data,
               chains = 3,
               iter = 500)
```

```
SAMPLING FOR MODEL 'anon_model' NOW (CHAIN 1).
Chain 1:
Chain 1: Gradient evaluation took 5.6e-05 seconds
Chain 1: 1000 transitions using 10 leapfrog steps per transition would take 0.56 seconds.
Chain 1: Adjust your expectations accordingly!
Chain 1:
Chain 1:
Chain 1: Iteration:   1 / 500 [  0%]  (Warmup)
Chain 1: Iteration:  50 / 500 [ 10%]  (Warmup)
Chain 1: Iteration: 100 / 500 [ 20%]  (Warmup)
Chain 1: Iteration: 150 / 500 [ 30%]  (Warmup)
Chain 1: Iteration: 200 / 500 [ 40%]  (Warmup)
Chain 1: Iteration: 250 / 500 [ 50%]  (Warmup)
Chain 1: Iteration: 251 / 500 [ 50%]  (Sampling)
Chain 1: Iteration: 300 / 500 [ 60%]  (Sampling)
Chain 1: Iteration: 350 / 500 [ 70%]  (Sampling)
Chain 1: Iteration: 400 / 500 [ 80%]  (Sampling)
Chain 1: Iteration: 450 / 500 [ 90%]  (Sampling)
Chain 1: Iteration: 500 / 500 [100%]  (Sampling)
Chain 1:
Chain 1:  Elapsed Time: 0.028 seconds (Warm-up)
Chain 1:                0.011 seconds (Sampling)
Chain 1:                0.039 seconds (Total)
Chain 1:

SAMPLING FOR MODEL 'anon_model' NOW (CHAIN 2).
Chain 2:
Chain 2: Gradient evaluation took 7e-06 seconds
Chain 2: 1000 transitions using 10 leapfrog steps per transition would take 0.07 seconds.
Chain 2: Adjust your expectations accordingly!
Chain 2:
Chain 2:
```

```
Chain 2: Iteration:   1 / 500 [  0%]  (Warmup)
Chain 2: Iteration:  50 / 500 [ 10%]  (Warmup)
Chain 2: Iteration: 100 / 500 [ 20%]  (Warmup)
Chain 2: Iteration: 150 / 500 [ 30%]  (Warmup)
Chain 2: Iteration: 200 / 500 [ 40%]  (Warmup)
Chain 2: Iteration: 250 / 500 [ 50%]  (Warmup)
Chain 2: Iteration: 251 / 500 [ 50%]  (Sampling)
Chain 2: Iteration: 300 / 500 [ 60%]  (Sampling)
Chain 2: Iteration: 350 / 500 [ 70%]  (Sampling)
Chain 2: Iteration: 400 / 500 [ 80%]  (Sampling)
Chain 2: Iteration: 450 / 500 [ 90%]  (Sampling)
Chain 2: Iteration: 500 / 500 [100%]  (Sampling)
Chain 2:
Chain 2:  Elapsed Time: 0.024 seconds (Warm-up)
Chain 2:                0.006 seconds (Sampling)
Chain 2:                0.03 seconds (Total)
Chain 2:


SAMPLING FOR MODEL 'anon_model' NOW (CHAIN 3).
Chain 3:
Chain 3: Gradient evaluation took 8e-06 seconds
Chain 3: 1000 transitions using 10 leapfrog steps per transition would take 0.08 seconds.
Chain 3: Adjust your expectations accordingly!
Chain 3:
Chain 3:
Chain 3: Iteration:   1 / 500 [  0%]  (Warmup)
Chain 3: Iteration:  50 / 500 [ 10%]  (Warmup)
Chain 3: Iteration: 100 / 500 [ 20%]  (Warmup)
Chain 3: Iteration: 150 / 500 [ 30%]  (Warmup)
Chain 3: Iteration: 200 / 500 [ 40%]  (Warmup)
Chain 3: Iteration: 250 / 500 [ 50%]  (Warmup)
Chain 3: Iteration: 251 / 500 [ 50%]  (Sampling)
Chain 3: Iteration: 300 / 500 [ 60%]  (Sampling)
Chain 3: Iteration: 350 / 500 [ 70%]  (Sampling)
Chain 3: Iteration: 400 / 500 [ 80%]  (Sampling)
Chain 3: Iteration: 450 / 500 [ 90%]  (Sampling)
Chain 3: Iteration: 500 / 500 [100%]  (Sampling)
Chain 3:
Chain 3:  Elapsed Time: 0.035 seconds (Warm-up)
Chain 3:                0.013 seconds (Sampling)
Chain 3:                0.048 seconds (Total)
Chain 3:
```

Look at the summary

```
fit
```

```
Inference for Stan model: anon_model.
3 chains, each with iter=500; warmup=250; thin=1;
post-warmup draws per chain=250, total post-warmup draws=750.

          mean se_mean   sd     2.5%      25%      50%      75%    97.5% n_eff
mu       86.74    0.04 0.93    84.79    86.12    86.76    87.42    88.34   475
sigma    20.38    0.03 0.70    19.04    19.89    20.37    20.86    21.87   575
lp__  -1525.73    0.05 0.93 -1528.06 -1526.10 -1525.45 -1525.06 -1524.78   391
      Rhat
mu    1.01
sigma 1.00
lp__  1.00

Samples were drawn using NUTS(diag_e) at Fri Feb 16 23:11:06 2024.
For each parameter, n_eff is a crude measure of effective sample size,
and Rhat is the potential scale reduction factor on split chains (at
convergence, Rhat=1).
```
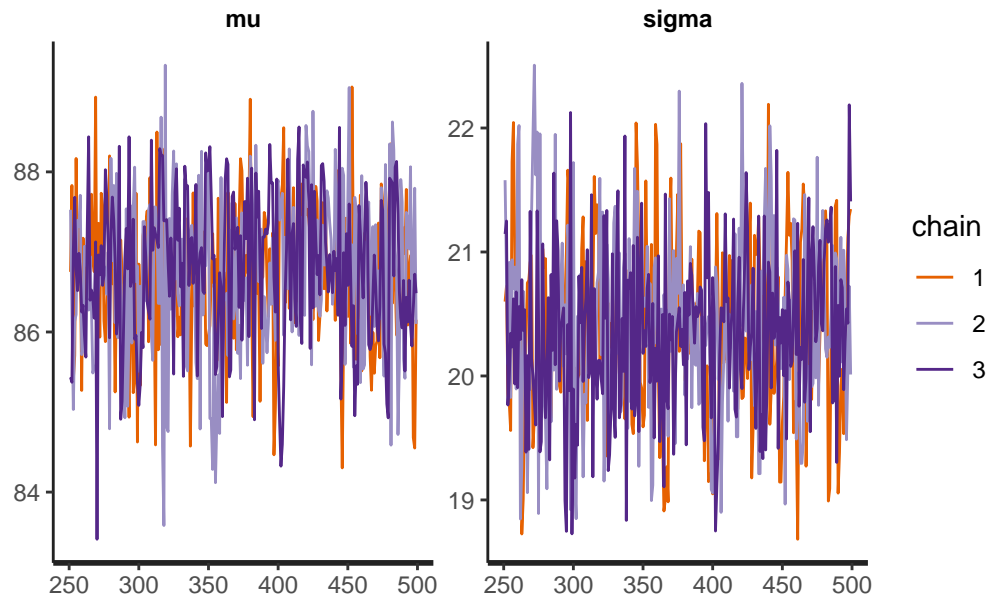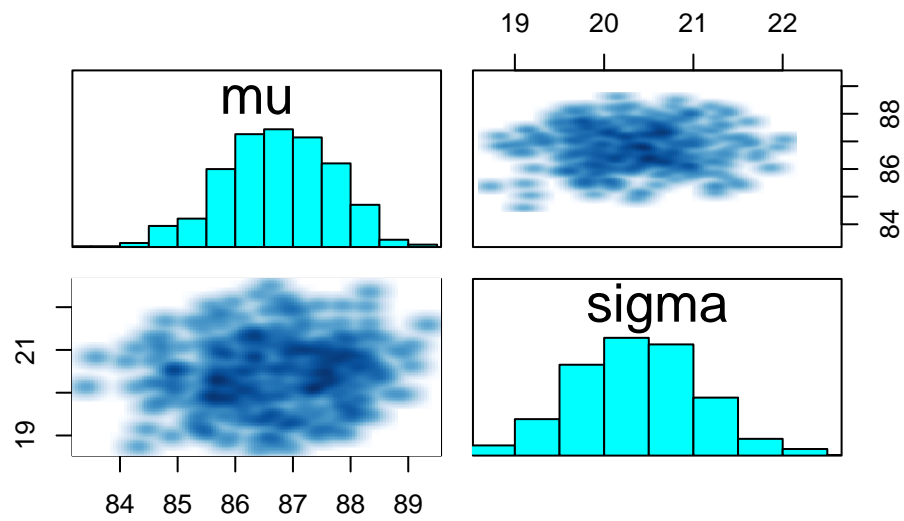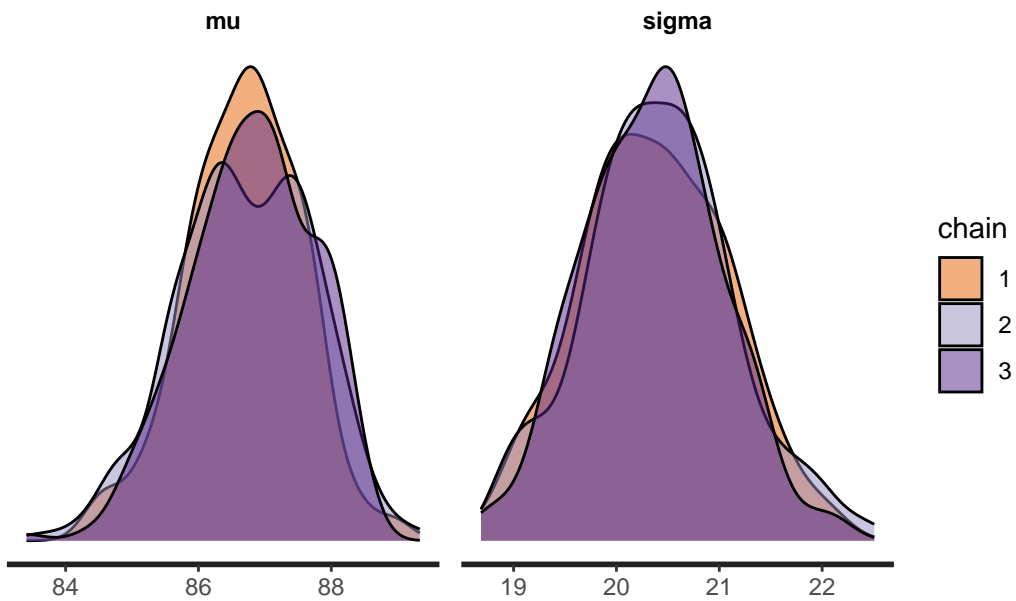
Traceplot

```
traceplot(fit)
```

All looks fine.

```
pairs(fit, pars = c("mu", "sigma"))
```

```
stan_dens(fit, separate_chains = TRUE)
```

## Understanding output

What does the model actually give us? A number of samples from the posteriors. To see this, we can use `extract` to get the samples.

```
post_samples <- extract(fit)
head(post_samples[["mu"]])
```

```
[1] 86.18439 87.74330 87.72298 86.66378 86.60305 85.27034
```

This is a list, and in this case, each element of the list has 4000 samples. E.g. quickly plot a histogram of mu

```
hist(post_samples[["mu"]])
```



**Histogram of post_samples[["mu"]]**

```
median(post_samples[["mu"]])
```

```
[1] 86.76287
```

```r
# 95% bayesian credible interval
quantile(post_samples[["mu"]], 0.025)
```

```
    2.5%
84.79048
```

```r
quantile(post_samples[["mu"]], 0.975)
```

```
   97.5%
88.33527
```

## Plot estimates

There are a bunch of packages, built-in functions that let you plot the estimates from the model, and I encourage you to explore these options (particularly in `bayesplot`, which we will most likely be using later on). I like using the `tidybayes` package, which allows us to easily get the posterior samples in a tidy format (e.g. using gather draws to get in long format). Once we have that, it's easy to just pipe and do ggplots as usual.

Get the posterior samples for mu and sigma in long format:

```r
dsamples <- fit  |>
  gather_draws(mu, sigma) # gather = long format
dsamples
```

```
# A tibble: 1,500 x 5
# Groups:   .variable [2]
   .chain .iteration .draw .variable .value
    <int>      <int> <int> <chr>      <dbl>
 1      1          1     1 mu          86.8
 2      1          2     2 mu          87.8
 3      1          3     3 mu          85.4
 4      1          4     4 mu          85.7
 5      1          5     5 mu          88.2
 6      1          6     6 mu          86.6
 7      1          7     7 mu          87.4
 8      1          8     8 mu          87.1
 9      1          9     9 mu          85.3
10      1         10    10 mu          87.2
# i 1,490 more rows
```

```
# wide format
fit |> spread_draws(mu, sigma)
```

```
# A tibble: 750 x 5
   .chain .iteration .draw    mu sigma
    <int>      <int> <int> <dbl> <dbl>
 1      1          1     1  86.8  20.6
 2      1          2     2  87.8  20.7
 3      1          3     3  85.4  19.8
 4      1          4     4  85.7  19.8
 5      1          5     5  88.2  19.6
 6      1          6     6  86.6  21.7
 7      1          7     7  87.4  22.0
 8      1          8     8  87.1  20.9
 9      1          9     9  85.3  20.5
10      1         10    10  87.2  20.5
# i 740 more rows
```

```
# quickly calculate the quantiles using

dsamples |>
  median_qi(.width = 0.8)
```
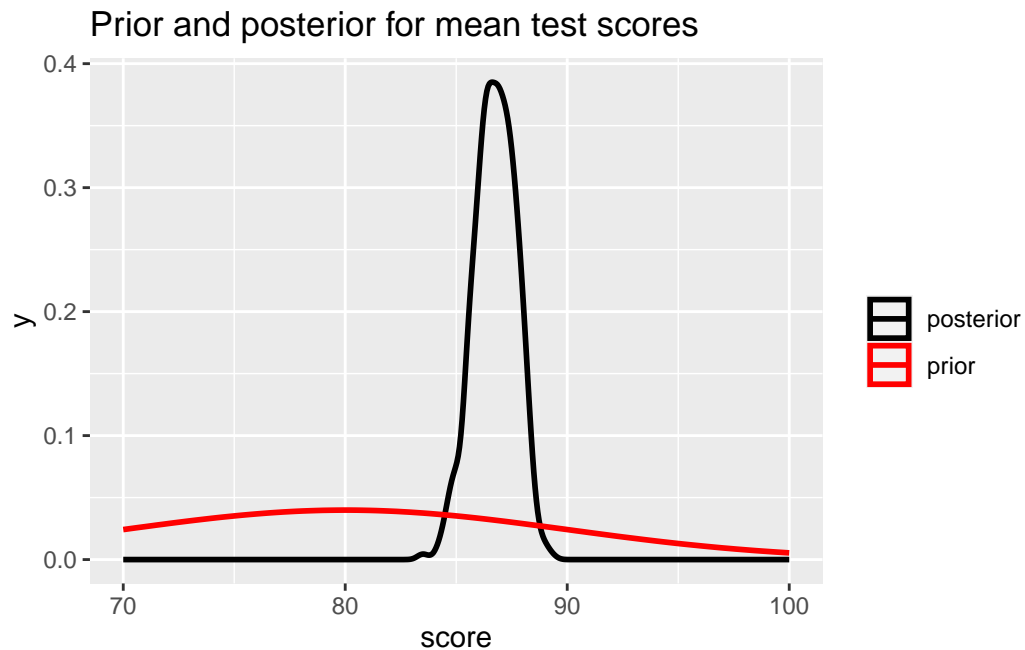
```
# A tibble: 2 x 7
  .variable .value .lower .upper .width .point .interval
  <chr>      <dbl>  <dbl>  <dbl>  <dbl> <chr>  <chr>
1 mu          86.8   85.6   87.9    0.8 median qi
2 sigma       20.4   19.5   21.3    0.8 median qi
```

Let's plot the density of the posterior samples for mu and add in the prior distribution

```
dsamples |>
  filter(.variable == "mu") |>
  ggplot(aes(.value, color = "posterior")) + geom_density(size = 1) +
  xlim(c(70, 100)) +
  stat_function(fun = dnorm,
        args = list(mean = mu0,
                    sd = sigma0),
        aes(colour = 'prior'), size = 1) +
  scale_color_manual(name = "", values = c("prior" = "red", "posterior" = "black")) +
```

```
ggtitle("Prior and posterior for mean test scores") +
xlab("score")
```

## Prior and posterior for mean test scores



## Question 2

Change the prior to be much more informative (by changing the standard deviation to be 0.1).
Rerun the model. Do the estimates change? Plot the prior and posterior densities.

```
mu0 <- 80
sigma0 <- 0.1
data <- list(y = y,
             N = length(y),
             mu0 = mu0,
             sigma0 = sigma0)
fit <- stan(file = "code/models/kids2.stan",
            data = data)
```

```
SAMPLING FOR MODEL 'anon_model' NOW (CHAIN 1).
Chain 1:
Chain 1: Gradient evaluation took 1e-05 seconds
```

```
Chain 1: 1000 transitions using 10 leapfrog steps per transition would take 0.1 seconds.
Chain 1: Adjust your expectations accordingly!
Chain 1:
Chain 1:
Chain 1: Iteration:    1 / 2000 [  0%]  (Warmup)
Chain 1: Iteration:  200 / 2000 [ 10%]  (Warmup)
Chain 1: Iteration:  400 / 2000 [ 20%]  (Warmup)
Chain 1: Iteration:  600 / 2000 [ 30%]  (Warmup)
Chain 1: Iteration:  800 / 2000 [ 40%]  (Warmup)
Chain 1: Iteration: 1000 / 2000 [ 50%]  (Warmup)
Chain 1: Iteration: 1001 / 2000 [ 50%]  (Sampling)
Chain 1: Iteration: 1200 / 2000 [ 60%]  (Sampling)
Chain 1: Iteration: 1400 / 2000 [ 70%]  (Sampling)
Chain 1: Iteration: 1600 / 2000 [ 80%]  (Sampling)
Chain 1: Iteration: 1800 / 2000 [ 90%]  (Sampling)
Chain 1: Iteration: 2000 / 2000 [100%]  (Sampling)
Chain 1:
Chain 1:  Elapsed Time: 0.035 seconds (Warm-up)
Chain 1:                0.031 seconds (Sampling)
Chain 1:                0.066 seconds (Total)
Chain 1:


SAMPLING FOR MODEL 'anon_model' NOW (CHAIN 2).
Chain 2:
Chain 2: Gradient evaluation took 7e-06 seconds
Chain 2: 1000 transitions using 10 leapfrog steps per transition would take 0.07 seconds.
Chain 2: Adjust your expectations accordingly!
Chain 2:
Chain 2:
Chain 2: Iteration:    1 / 2000 [  0%]  (Warmup)
Chain 2: Iteration:  200 / 2000 [ 10%]  (Warmup)
Chain 2: Iteration:  400 / 2000 [ 20%]  (Warmup)
Chain 2: Iteration:  600 / 2000 [ 30%]  (Warmup)
Chain 2: Iteration:  800 / 2000 [ 40%]  (Warmup)
Chain 2: Iteration: 1000 / 2000 [ 50%]  (Warmup)
Chain 2: Iteration: 1001 / 2000 [ 50%]  (Sampling)
Chain 2: Iteration: 1200 / 2000 [ 60%]  (Sampling)
Chain 2: Iteration: 1400 / 2000 [ 70%]  (Sampling)
Chain 2: Iteration: 1600 / 2000 [ 80%]  (Sampling)
Chain 2: Iteration: 1800 / 2000 [ 90%]  (Sampling)
Chain 2: Iteration: 2000 / 2000 [100%]  (Sampling)
Chain 2:
Chain 2:  Elapsed Time: 0.035 seconds (Warm-up)
```

```
Chain 2:                    0.031 seconds (Sampling)
Chain 2:                    0.066 seconds (Total)
Chain 2:


SAMPLING FOR MODEL 'anon_model' NOW (CHAIN 3).
Chain 3:
Chain 3: Gradient evaluation took 7e-06 seconds
Chain 3: 1000 transitions using 10 leapfrog steps per transition would take 0.07 seconds.
Chain 3: Adjust your expectations accordingly!
Chain 3:
Chain 3:
Chain 3: Iteration:    1 / 2000 [  0%]  (Warmup)
Chain 3: Iteration:  200 / 2000 [ 10%]  (Warmup)
Chain 3: Iteration:  400 / 2000 [ 20%]  (Warmup)
Chain 3: Iteration:  600 / 2000 [ 30%]  (Warmup)
Chain 3: Iteration:  800 / 2000 [ 40%]  (Warmup)
Chain 3: Iteration: 1000 / 2000 [ 50%]  (Warmup)
Chain 3: Iteration: 1001 / 2000 [ 50%]  (Sampling)
Chain 3: Iteration: 1200 / 2000 [ 60%]  (Sampling)
Chain 3: Iteration: 1400 / 2000 [ 70%]  (Sampling)
Chain 3: Iteration: 1600 / 2000 [ 80%]  (Sampling)
Chain 3: Iteration: 1800 / 2000 [ 90%]  (Sampling)
Chain 3: Iteration: 2000 / 2000 [100%]  (Sampling)
Chain 3:
Chain 3:  Elapsed Time: 0.033 seconds (Warm-up)
Chain 3:                0.043 seconds (Sampling)
Chain 3:                0.076 seconds (Total)
Chain 3:


SAMPLING FOR MODEL 'anon_model' NOW (CHAIN 4).
Chain 4:
Chain 4: Gradient evaluation took 7e-06 seconds
Chain 4: 1000 transitions using 10 leapfrog steps per transition would take 0.07 seconds.
Chain 4: Adjust your expectations accordingly!
Chain 4:
Chain 4:
Chain 4: Iteration:    1 / 2000 [  0%]  (Warmup)
Chain 4: Iteration:  200 / 2000 [ 10%]  (Warmup)
Chain 4: Iteration:  400 / 2000 [ 20%]  (Warmup)
Chain 4: Iteration:  600 / 2000 [ 30%]  (Warmup)
Chain 4: Iteration:  800 / 2000 [ 40%]  (Warmup)
Chain 4: Iteration: 1000 / 2000 [ 50%]  (Warmup)
Chain 4: Iteration: 1001 / 2000 [ 50%]  (Sampling)
```

```
Chain 4: Iteration: 1200 / 2000 [ 60%]  (Sampling)
Chain 4: Iteration: 1400 / 2000 [ 70%]  (Sampling)
Chain 4: Iteration: 1600 / 2000 [ 80%]  (Sampling)
Chain 4: Iteration: 1800 / 2000 [ 90%]  (Sampling)
Chain 4: Iteration: 2000 / 2000 [100%]  (Sampling)
Chain 4:
Chain 4:  Elapsed Time: 0.045 seconds (Warm-up)
Chain 4:                0.029 seconds (Sampling)
Chain 4:                0.074 seconds (Total)
Chain 4:
```

```
fit
```

```
Inference for Stan model: anon_model.
4 chains, each with iter=2000; warmup=1000; thin=1;
post-warmup draws per chain=1000, total post-warmup draws=4000.

          mean se_mean   sd     2.5%       25%       50%       75%     97.5% n_eff
mu        80.06    0.00 0.10    79.86     80.00     80.06     80.13     80.26  3621
sigma     21.43    0.01 0.70    20.09     20.95     21.40     21.90     22.82  3623
lp__   -1548.34    0.02 0.98 -1550.87 -1548.72 -1548.05 -1547.65 -1547.40  1931
       Rhat
mu        1
sigma     1
lp__      1

Samples were drawn using NUTS(diag_e) at Fri Feb 16 23:11:11 2024.
For each parameter, n_eff is a crude measure of effective sample size,
and Rhat is the potential scale reduction factor on split chains (at
convergence, Rhat=1).
```

Yes, both estimates of mu and sigma increase.

```
dsamples <- fit  |>
  gather_draws(mu, sigma)
```
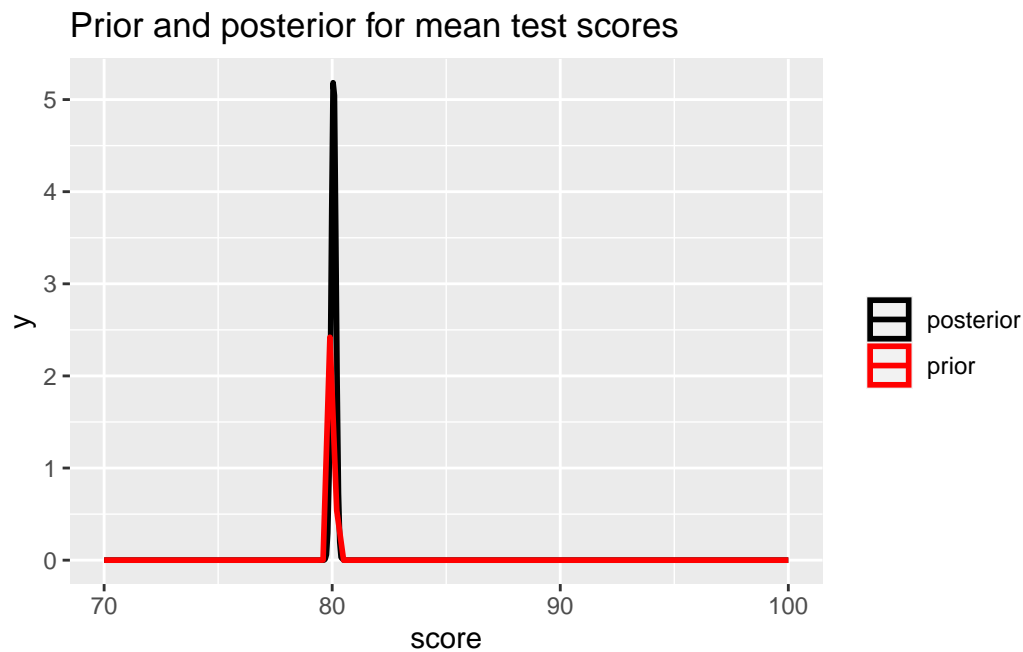
```
dsamples |>
  filter(.variable == "mu") |>
  ggplot(aes(.value, color = "posterior")) + geom_density(size = 1) +
  xlim(c(70, 100)) +
  stat_function(fun = dnorm,
```

```
        args = list(mean = mu0,
                    sd = sigma0),
        aes(colour = 'prior'), size = 1) +
  scale_color_manual(name = "", values = c("prior" = "red", "posterior" = "black")) +
  ggtitle("Prior and posterior for mean test scores") +
  xlab("score")
```

## Prior and posterior for mean test scores



## Adding covariates

Now let's see how kid's test scores are related to mother's education. We want to run the simple linear regression

$$Score = \alpha + \beta X$$

where $X = 1$ if the mother finished high school and zero otherwise.

`kid3.stan` has the stan model to do this. Notice now we have some inputs related to the design matrix $X$ and the number of covariates (in this case, it's just 1).

Let's get the data we need and run the model.

```r
X <- as.matrix(kidiq$mom_hs, ncol = 1) # force this to be a matrix
K <- 1

data <- list(y = y, N = length(y),
             X =X, K = K)
fit2 <- stan(file = here("code/models/kids3.stan"),
             data = data,
             iter = 1000)
```

```
SAMPLING FOR MODEL 'anon_model' NOW (CHAIN 1).
Chain 1:
Chain 1: Gradient evaluation took 0.000115 seconds
Chain 1: 1000 transitions using 10 leapfrog steps per transition would take 1.15 seconds.
Chain 1: Adjust your expectations accordingly!
Chain 1:
Chain 1:
Chain 1: Iteration:    1 / 1000 [  0%]  (Warmup)
Chain 1: Iteration: 100 / 1000 [ 10%]  (Warmup)
Chain 1: Iteration: 200 / 1000 [ 20%]  (Warmup)
Chain 1: Iteration: 300 / 1000 [ 30%]  (Warmup)
Chain 1: Iteration: 400 / 1000 [ 40%]  (Warmup)
Chain 1: Iteration: 500 / 1000 [ 50%]  (Warmup)
Chain 1: Iteration: 501 / 1000 [ 50%]  (Sampling)
Chain 1: Iteration: 600 / 1000 [ 60%]  (Sampling)
Chain 1: Iteration: 700 / 1000 [ 70%]  (Sampling)
Chain 1: Iteration: 800 / 1000 [ 80%]  (Sampling)
Chain 1: Iteration: 900 / 1000 [ 90%]  (Sampling)
Chain 1: Iteration: 1000 / 1000 [100%]  (Sampling)
Chain 1:
Chain 1:  Elapsed Time: 0.282 seconds (Warm-up)
Chain 1:                0.126 seconds (Sampling)
Chain 1:                0.408 seconds (Total)
Chain 1:

SAMPLING FOR MODEL 'anon_model' NOW (CHAIN 2).
Chain 2:
Chain 2: Gradient evaluation took 2.8e-05 seconds
Chain 2: 1000 transitions using 10 leapfrog steps per transition would take 0.28 seconds.
Chain 2: Adjust your expectations accordingly!
Chain 2:
Chain 2:
```

```
Chain 2: Iteration:    1 / 1000 [  0%]   (Warmup)
Chain 2: Iteration: 100 / 1000 [ 10%]   (Warmup)
Chain 2: Iteration: 200 / 1000 [ 20%]   (Warmup)
Chain 2: Iteration: 300 / 1000 [ 30%]   (Warmup)
Chain 2: Iteration: 400 / 1000 [ 40%]   (Warmup)
Chain 2: Iteration: 500 / 1000 [ 50%]   (Warmup)
Chain 2: Iteration: 501 / 1000 [ 50%]   (Sampling)
Chain 2: Iteration: 600 / 1000 [ 60%]   (Sampling)
Chain 2: Iteration: 700 / 1000 [ 70%]   (Sampling)
Chain 2: Iteration: 800 / 1000 [ 80%]   (Sampling)
Chain 2: Iteration: 900 / 1000 [ 90%]   (Sampling)
Chain 2: Iteration: 1000 / 1000 [100%]   (Sampling)
Chain 2:
Chain 2:   Elapsed Time: 0.23 seconds (Warm-up)
Chain 2:                 0.158 seconds (Sampling)
Chain 2:                 0.388 seconds (Total)
Chain 2:


SAMPLING FOR MODEL 'anon_model' NOW (CHAIN 3).
Chain 3:
Chain 3: Gradient evaluation took 3.9e-05 seconds
Chain 3: 1000 transitions using 10 leapfrog steps per transition would take 0.39 seconds.
Chain 3: Adjust your expectations accordingly!
Chain 3:
Chain 3:
Chain 3: Iteration:    1 / 1000 [  0%]   (Warmup)
Chain 3: Iteration: 100 / 1000 [ 10%]   (Warmup)
Chain 3: Iteration: 200 / 1000 [ 20%]   (Warmup)
Chain 3: Iteration: 300 / 1000 [ 30%]   (Warmup)
Chain 3: Iteration: 400 / 1000 [ 40%]   (Warmup)
Chain 3: Iteration: 500 / 1000 [ 50%]   (Warmup)
Chain 3: Iteration: 501 / 1000 [ 50%]   (Sampling)
Chain 3: Iteration: 600 / 1000 [ 60%]   (Sampling)
Chain 3: Iteration: 700 / 1000 [ 70%]   (Sampling)
Chain 3: Iteration: 800 / 1000 [ 80%]   (Sampling)
Chain 3: Iteration: 900 / 1000 [ 90%]   (Sampling)
Chain 3: Iteration: 1000 / 1000 [100%]   (Sampling)
Chain 3:
Chain 3:   Elapsed Time: 0.221 seconds (Warm-up)
Chain 3:                 0.142 seconds (Sampling)
Chain 3:                 0.363 seconds (Total)
Chain 3:
```

```
SAMPLING FOR MODEL 'anon_model' NOW (CHAIN 4).
Chain 4:
Chain 4: Gradient evaluation took 3.2e-05 seconds
Chain 4: 1000 transitions using 10 leapfrog steps per transition would take 0.32 seconds.
Chain 4: Adjust your expectations accordingly!
Chain 4:
Chain 4:
Chain 4: Iteration:    1 / 1000 [  0%]  (Warmup)
Chain 4: Iteration: 100 / 1000 [ 10%]  (Warmup)
Chain 4: Iteration: 200 / 1000 [ 20%]  (Warmup)
Chain 4: Iteration: 300 / 1000 [ 30%]  (Warmup)
Chain 4: Iteration: 400 / 1000 [ 40%]  (Warmup)
Chain 4: Iteration: 500 / 1000 [ 50%]  (Warmup)
Chain 4: Iteration: 501 / 1000 [ 50%]  (Sampling)
Chain 4: Iteration: 600 / 1000 [ 60%]  (Sampling)
Chain 4: Iteration: 700 / 1000 [ 70%]  (Sampling)
Chain 4: Iteration: 800 / 1000 [ 80%]  (Sampling)
Chain 4: Iteration: 900 / 1000 [ 90%]  (Sampling)
Chain 4: Iteration: 1000 / 1000 [100%]  (Sampling)
Chain 4:
Chain 4:  Elapsed Time: 0.283 seconds (Warm-up)
Chain 4:                0.162 seconds (Sampling)
Chain 4:                0.445 seconds (Total)
Chain 4:
```

## Question 3

a) Confirm that the estimates of the intercept and slope are comparable to results from `lm()`

b) Do a `pairs` plot to investigate the joint sample distributions of the slope and intercept. Comment briefly on what you see. Is this potentially a problem?

```
summary(fit2)$summary
```

|  | mean | se_mean | sd | 2.5% | 25% | 50% |
|---|---|---|---|---|---|---|
| alpha | 77.89094 | 0.06763012 | 2.0036297 | 73.908666 | 76.541717 | 77.86273 |
| beta[1] | 11.32511 | 0.07448299 | 2.2285360 | 6.874445 | 9.941839 | 11.30399 |
| sigma | 19.80873 | 0.02278819 | 0.6959131 | 18.423953 | 19.350433 | 19.80171 |
| lp__ | -1514.39317 | 0.04402574 | 1.2785897 | -1517.733700 | -1514.944376 | -1514.04596 |

|  | 75% | 97.5% | n_eff | Rhat |
|---|---|---|---|---|
| alpha | 79.19361 | 82.00795 | 877.7172 | 1.001857 |

```
beta[1]      12.75657     15.59413 895.2103 1.002523
sigma        20.24819     21.21316 932.5894 1.003807
lp__       -1513.48358 -1512.97367 843.4302 1.004699
```
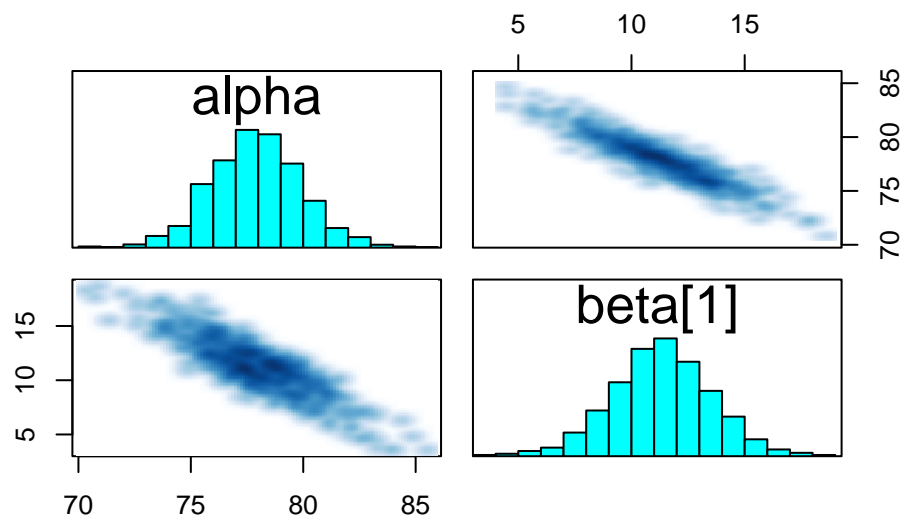
```r
lm(y~kidiq$mom_hs)
```

```
Call:
lm(formula = y ~ kidiq$mom_hs)

Coefficients:
 (Intercept)  kidiq$mom_hs
       77.55         11.77
```

They have quiet similar estimates.

```r
pairs(fit2, pars = c("alpha", "beta[1]"))
```
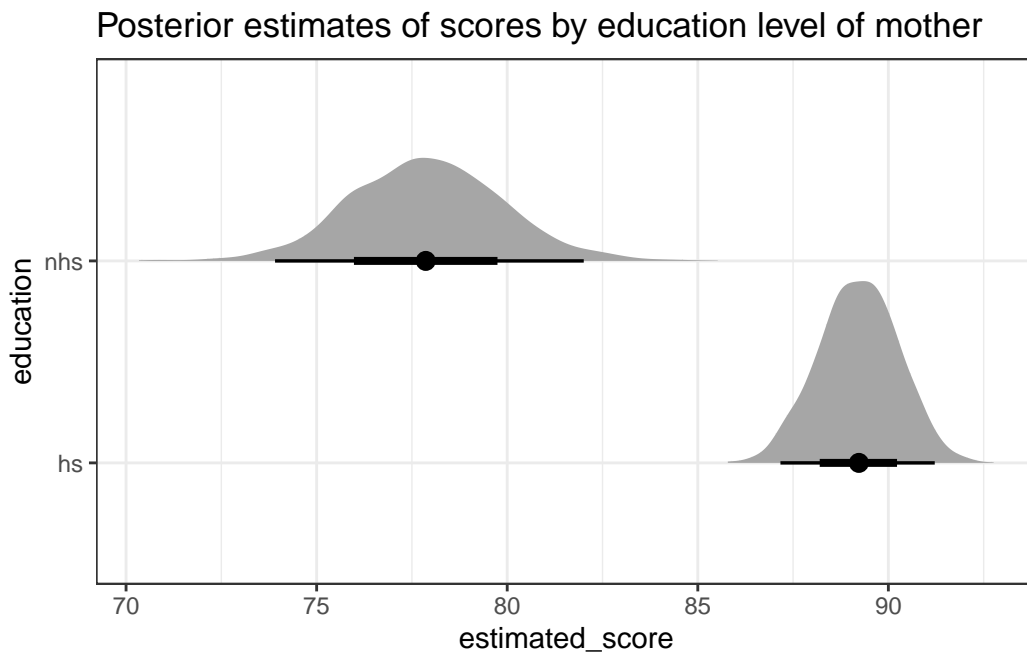


The coefficients are strongly correlated in a negative direction. This can be problematic because you may not obtain a wider range of samples.

**Plotting results**

It might be nice to plot the posterior samples of the estimates for the non-high-school and high-school mothered kids. Here's some code that does this: notice the `beta[condition]` syntax. Also notice I'm using `spread_draws`, because it's easier to calculate the estimated effects in wide format

```
fit2 |>
  spread_draws(alpha, beta[k], sigma) |>
    mutate(nhs = alpha, # no high school is just the intercept
           hs = alpha + beta) |>
  select(nhs, hs) |>
  pivot_longer(nhs:hs, names_to = "education", values_to = "estimated_score") |>
  ggplot(aes(y = education, x = estimated_score)) +
  stat_halfeye() +
  theme_bw() +
  ggtitle("Posterior estimates of scores by education level of mother")
```



Posterior estimates of scores by education level of mother

**Question 4**

Add in mother's IQ as a covariate and rerun the model. Please mean center the covariate before putting it into the model. Interpret the coefficient on the (centered) mum's IQ.

```r
X <- cbind(kidiq$mom_hs, kidiq$mom_iq - mean(kidiq$mom_iq))
data <- list(y = y,
             N = length(y),
             K = 2,
             X = as.matrix(X))


fit2 <- stan(file = "code/models/kids3.stan",
             data = data,
             iter = 1000)
```

```
SAMPLING FOR MODEL 'anon_model' NOW (CHAIN 1).
Chain 1:
Chain 1: Gradient evaluation took 3.4e-05 seconds
Chain 1: 1000 transitions using 10 leapfrog steps per transition would take 0.34 seconds.
Chain 1: Adjust your expectations accordingly!
Chain 1:
Chain 1:
Chain 1: Iteration:    1 / 1000 [  0%]  (Warmup)
Chain 1: Iteration: 100 / 1000 [ 10%]  (Warmup)
Chain 1: Iteration: 200 / 1000 [ 20%]  (Warmup)
Chain 1: Iteration: 300 / 1000 [ 30%]  (Warmup)
Chain 1: Iteration: 400 / 1000 [ 40%]  (Warmup)
Chain 1: Iteration: 500 / 1000 [ 50%]  (Warmup)
Chain 1: Iteration: 501 / 1000 [ 50%]  (Sampling)
Chain 1: Iteration: 600 / 1000 [ 60%]  (Sampling)
Chain 1: Iteration: 700 / 1000 [ 70%]  (Sampling)
Chain 1: Iteration: 800 / 1000 [ 80%]  (Sampling)
Chain 1: Iteration: 900 / 1000 [ 90%]  (Sampling)
Chain 1: Iteration: 1000 / 1000 [100%]  (Sampling)
Chain 1:
Chain 1:  Elapsed Time: 0.209 seconds (Warm-up)
Chain 1:                0.155 seconds (Sampling)
Chain 1:                0.364 seconds (Total)
Chain 1:

SAMPLING FOR MODEL 'anon_model' NOW (CHAIN 2).
Chain 2:
Chain 2: Gradient evaluation took 3e-05 seconds
Chain 2: 1000 transitions using 10 leapfrog steps per transition would take 0.3 seconds.
Chain 2: Adjust your expectations accordingly!
```

```
Chain 2:
Chain 2:
Chain 2: Iteration:    1 / 1000 [  0%]  (Warmup)
Chain 2: Iteration: 100 / 1000 [ 10%]  (Warmup)
Chain 2: Iteration: 200 / 1000 [ 20%]  (Warmup)
Chain 2: Iteration: 300 / 1000 [ 30%]  (Warmup)
Chain 2: Iteration: 400 / 1000 [ 40%]  (Warmup)
Chain 2: Iteration: 500 / 1000 [ 50%]  (Warmup)
Chain 2: Iteration: 501 / 1000 [ 50%]  (Sampling)
Chain 2: Iteration: 600 / 1000 [ 60%]  (Sampling)
Chain 2: Iteration: 700 / 1000 [ 70%]  (Sampling)
Chain 2: Iteration: 800 / 1000 [ 80%]  (Sampling)
Chain 2: Iteration: 900 / 1000 [ 90%]  (Sampling)
Chain 2: Iteration: 1000 / 1000 [100%]  (Sampling)
Chain 2:
Chain 2:  Elapsed Time: 0.186 seconds (Warm-up)
Chain 2:                0.146 seconds (Sampling)
Chain 2:                0.332 seconds (Total)
Chain 2:


SAMPLING FOR MODEL 'anon_model' NOW (CHAIN 3).
Chain 3:
Chain 3: Gradient evaluation took 2.9e-05 seconds
Chain 3: 1000 transitions using 10 leapfrog steps per transition would take 0.29 seconds.
Chain 3: Adjust your expectations accordingly!
Chain 3:
Chain 3:
Chain 3: Iteration:    1 / 1000 [  0%]  (Warmup)
Chain 3: Iteration: 100 / 1000 [ 10%]  (Warmup)
Chain 3: Iteration: 200 / 1000 [ 20%]  (Warmup)
Chain 3: Iteration: 300 / 1000 [ 30%]  (Warmup)
Chain 3: Iteration: 400 / 1000 [ 40%]  (Warmup)
Chain 3: Iteration: 500 / 1000 [ 50%]  (Warmup)
Chain 3: Iteration: 501 / 1000 [ 50%]  (Sampling)
Chain 3: Iteration: 600 / 1000 [ 60%]  (Sampling)
Chain 3: Iteration: 700 / 1000 [ 70%]  (Sampling)
Chain 3: Iteration: 800 / 1000 [ 80%]  (Sampling)
Chain 3: Iteration: 900 / 1000 [ 90%]  (Sampling)
Chain 3: Iteration: 1000 / 1000 [100%]  (Sampling)
Chain 3:
Chain 3:  Elapsed Time: 0.239 seconds (Warm-up)
Chain 3:                0.178 seconds (Sampling)
Chain 3:                0.417 seconds (Total)
```

```
Chain 3:

SAMPLING FOR MODEL 'anon_model' NOW (CHAIN 4).
Chain 4:
Chain 4: Gradient evaluation took 2.8e-05 seconds
Chain 4: 1000 transitions using 10 leapfrog steps per transition would take 0.28 seconds.
Chain 4: Adjust your expectations accordingly!
Chain 4:
Chain 4:
Chain 4: Iteration:    1 / 1000 [  0%]  (Warmup)
Chain 4: Iteration: 100 / 1000 [ 10%]  (Warmup)
Chain 4: Iteration: 200 / 1000 [ 20%]  (Warmup)
Chain 4: Iteration: 300 / 1000 [ 30%]  (Warmup)
Chain 4: Iteration: 400 / 1000 [ 40%]  (Warmup)
Chain 4: Iteration: 500 / 1000 [ 50%]  (Warmup)
Chain 4: Iteration: 501 / 1000 [ 50%]  (Sampling)
Chain 4: Iteration: 600 / 1000 [ 60%]  (Sampling)
Chain 4: Iteration: 700 / 1000 [ 70%]  (Sampling)
Chain 4: Iteration: 800 / 1000 [ 80%]  (Sampling)
Chain 4: Iteration: 900 / 1000 [ 90%]  (Sampling)
Chain 4: Iteration: 1000 / 1000 [100%]  (Sampling)
Chain 4:
Chain 4:  Elapsed Time: 0.236 seconds (Warm-up)
Chain 4:                0.184 seconds (Sampling)
Chain 4:                0.42 seconds (Total)
Chain 4:
```

```
  summary(fit2)$summary
```

```
                 mean       se_mean          sd         2.5%            25%
alpha      82.2919818 0.059566666 1.86435786    78.6139825    81.0701754
beta[1]     5.7670277 0.066170346 2.12105181     1.6669809     4.3525089
beta[2]     0.5639563 0.001546756 0.06168746     0.4443275     0.5200515
sigma      18.1134371 0.016450506 0.62758504    16.9154020    17.6927044
lp__    -1474.4708856 0.055201175 1.38213012 -1477.9542195 -1475.1844488
                 50%           75%         97.5%       n_eff      Rhat
alpha      82.2563264    83.5391708    86.0232300   979.6072 0.9991167
beta[1]     5.7713981     7.2121068     9.8449108  1027.4856 0.9989629
beta[2]     0.5639007     0.6078943     0.6838726  1590.5606 1.0029256
sigma      18.0844110    18.5237468    19.3902244  1455.4143 1.0030863
lp__    -1474.1831335 -1473.4367441 -1472.6757851   626.9043 0.9997542
```

When the centered mom's IQ increases by one unit, the expected kid's score increase by 0.5657213 with holding the other variables as constant.

## Question 5

Confirm the results from Stan agree with `lm()`

```
lm(y ~ X[,1] +  X[,2])
```

```
Call:
lm(formula = y ~ X[, 1] + X[, 2])

Coefficients:
(Intercept)        X[, 1]        X[, 2]
    82.1221        5.9501        0.5639
```

The coefficients are not exactly the same, but it is still consistent.
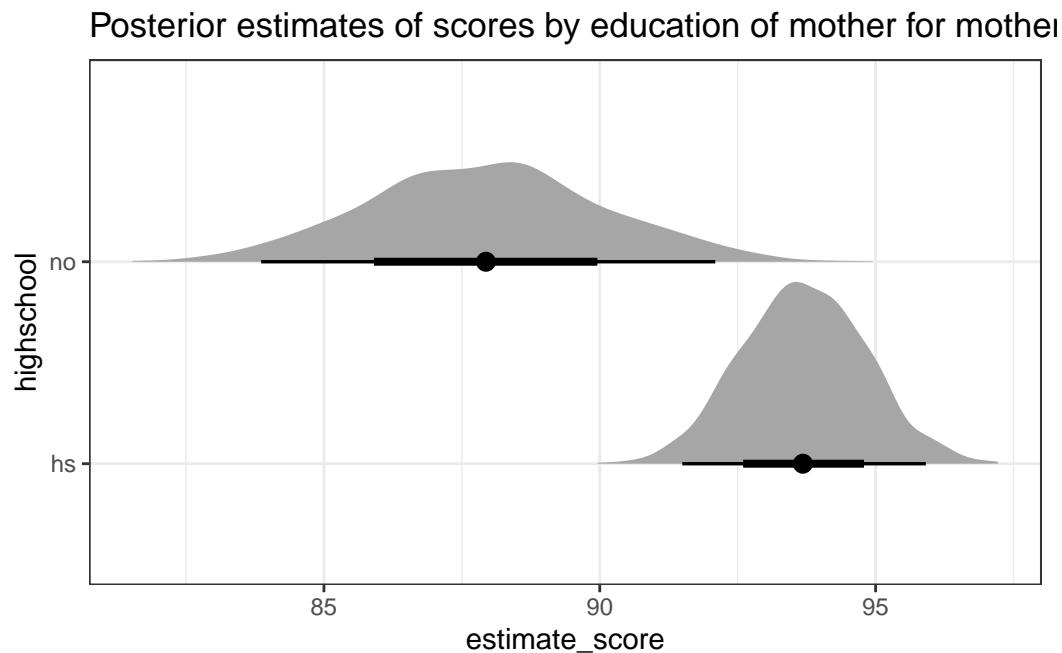
## Question 6

Plot the posterior estimates of scores by education of mother for mothers who have an IQ of 110.

```
posterior <- fit2 |>
  spread_draws(alpha, beta[k], sigma) |>
  mutate(mom_iq_adj = 110 - mean(kidiq$mom_iq)) |>
  mutate(estimate_no = alpha + mom_iq_adj * beta[2],
         estimate_hs = alpha + mom_iq_adj * beta[2] + beta[1]) |>
  select(estimate_no, estimate_hs) |>
  pivot_longer(estimate_no:estimate_hs, names_to = "highschool", values_to = "estimate_sco
```

```
fit2 |>
  spread_draws(alpha, beta[k], sigma) |>
  pivot_wider(names_from = k, values_from = beta,  names_glue = "beta{k}") |>
  mutate(no = alpha + beta2 * (110 - mean(kidiq$mom_iq)),
         hs = alpha + beta1 + beta2 * (110 - mean(kidiq$mom_iq))) |>
  select(no, hs) |>
  pivot_longer(no:hs, names_to = "highschool", values_to = "estimate_score") |>
  ggplot(aes(y = highschool, x = estimate_score)) +
```

```
stat_halfeye() +
theme_bw() +
ggtitle("Posterior estimates of scores by education of mother for mothers who have an IQ
```

## Posterior estimates of scores by education of mother for mother



### Question 7

Generate and plot (as a histogram) samples from the posterior predictive distribution for a
new kid with a mother who graduated high school and has an IQ of 95.

```
sample <- extract(fit2)
```

```
prediction_mean <- sample$alpha + sample$beta[,1] + (95 - mean(kidiq$mom_iq))*sample$beta[
y_pred <- tibble(y_pred = rnorm(length(sample$sigma),
                                mean = prediction_mean,
                                sd = sample$sigma))
```

```
ggplot(y_pred, aes(y_pred)) +
  geom_histogram(fill = "lightgreen", col = "black") +
  ggtitle("Distribution of Predicted Scores with a mother who graduated high school and ha
```

Distribution of Predicted Scores with a mother who graduated I