

# Week 5: Bayesian linear regression and introduction to Stan

Hainan Xu

15/02/24

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

```
kidiq <- readRDS("~/Documents/2024/STA2201-Methods-of-Applied-Statistics/labs/kidiq.RDS")
```

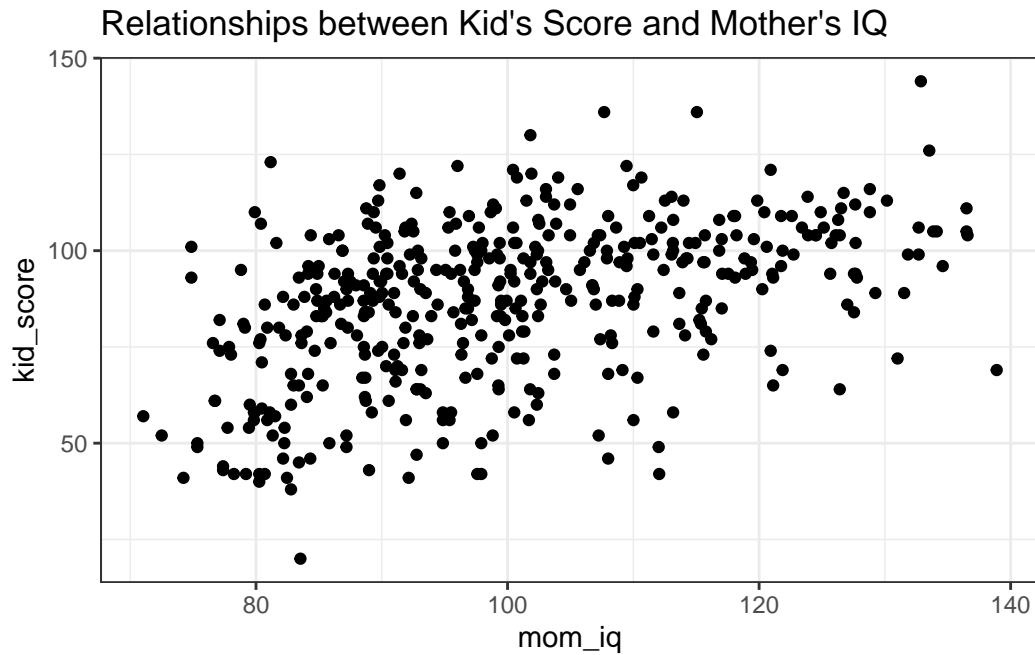
## Question 1

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

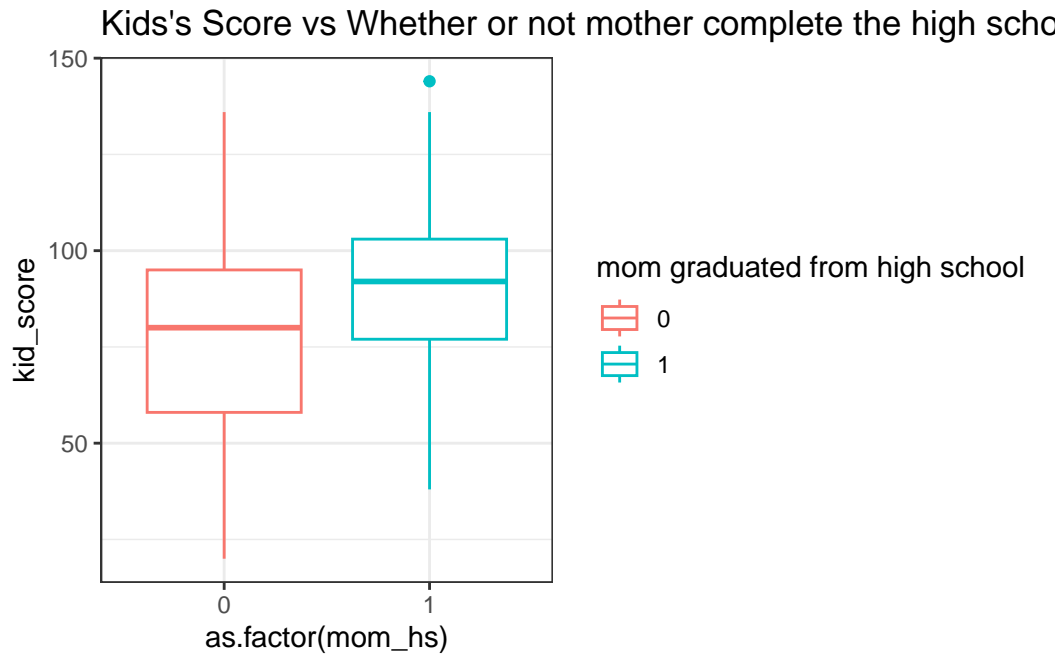
The scatterplot below shows the relationships between kid's score and mom's IQ, we observe that the kid's score and mom's IQ seems to have a positive relationship.

```
#relation ship between kids' score and mother's IQ
ggplot(data=kidiq)+geom_point(aes(x=mom_iq,y=kid_score),fill="lightblue")+theme_bw()+ggtitle
```



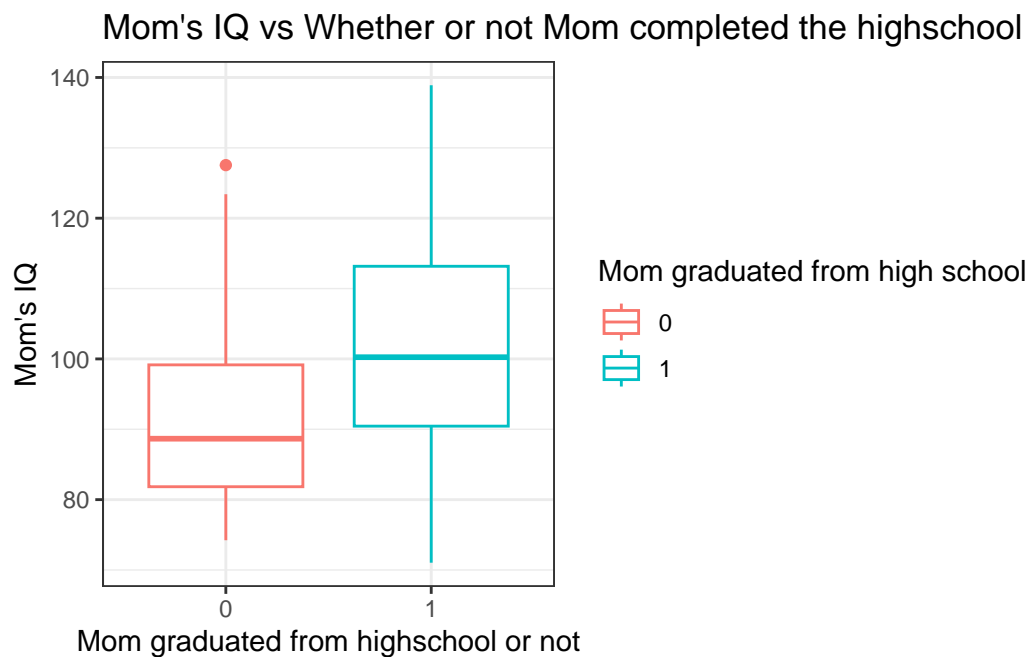
The boxplot below shows the distribution of kids' score categorized by whether or not the mom is graduated from highschool. We observe that the mean of kid's IQ whose graduated from highschool is higher than that of not graduated from highschool.

```
# relationships between kids'score and whether or not mother complete the high school  
ggplot(data=kidiq)+geom_boxplot(aes(x=as.factor(mom_hs),y=kid_score,color=as.factor(mom_hs)))
```



The boxplot below shows the distribution of mom's iq catergrized by whether or not the shei s graduated from highschool. We observe that the mean of mom's IQ who graduated from highschool is higher than that of not graduated from highschool. Interestingly, fomr the given data, we observe that the IQ of the highschool graduated mom can still be lower than those who did not graduate from highschool.

```
#relationships between mom's IQ and whether or not they completed high school
ggplot(data=kidiq)+geom_boxplot(aes(x=as.factor(mom_hs),y=mom_iq,color=as.factor(mom_hs)))
```



## Original Model

In order to compare the density plot as required in question2, we keep original model fitted with  $\sigma_0=0.1$ .

```
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)
```

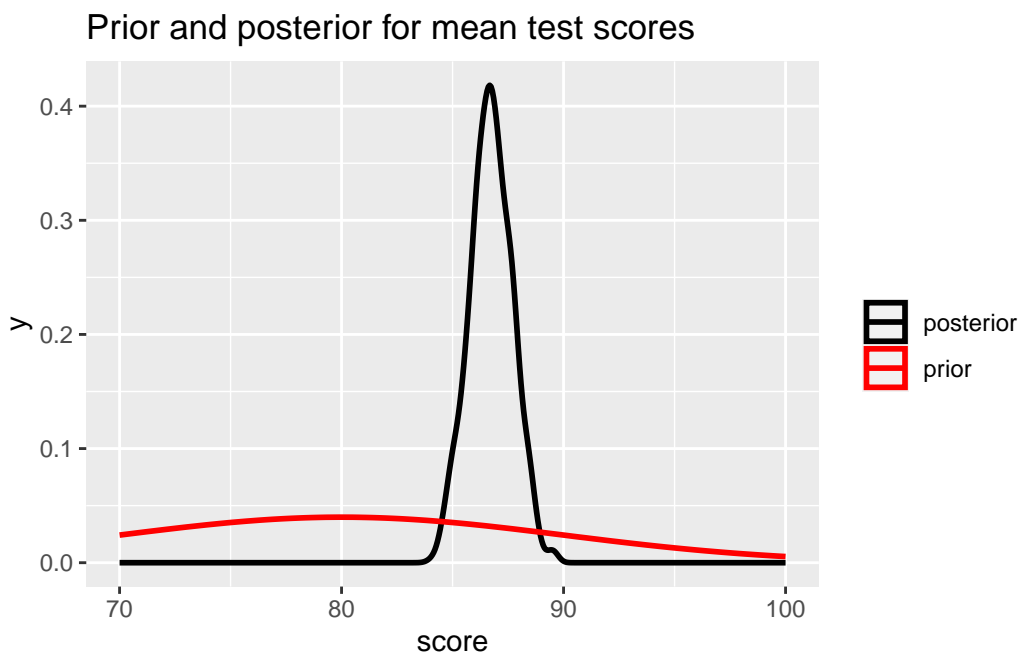
Now we can run the model:

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

```
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 1 mu      85.9
2      1      2      2 2 mu      87.6
3      1      3      3 3 mu      87.8
4      1      4      4 4 mu      89.6
5      1      5      5 5 mu      87.9
6      1      6      6 6 mu      86.6
7      1      7      7 7 mu      86.6
8      1      8      8 8 mu      85.8
9      1      9      9 9 mu      86.5
10     1     10     10 10 mu      87.0
# ... with 1,490 more rows
```

```
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")
```



## 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.

The estimates changed after changing the standard deviation as follows.  $\mu$  decreased from 86 to around 80.  $\sigma$  increased from 20 to 21.

```
y <- kidiq$kid_score
mu0 <- 80
sigma0 <- 0.1

# named list to input for stan function
data <- list(y = y,
             N = length(y),
             mu0 = mu0,
             sigma0 = sigma0)

#check the model fit
fit <- stan(file = here("~/Documents/2024/STA2201-Methods-of-Applied-Statistics/labs/kids2"),
            data = data,
            chains = 3,
```

```
iter = 500)
```

```
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.

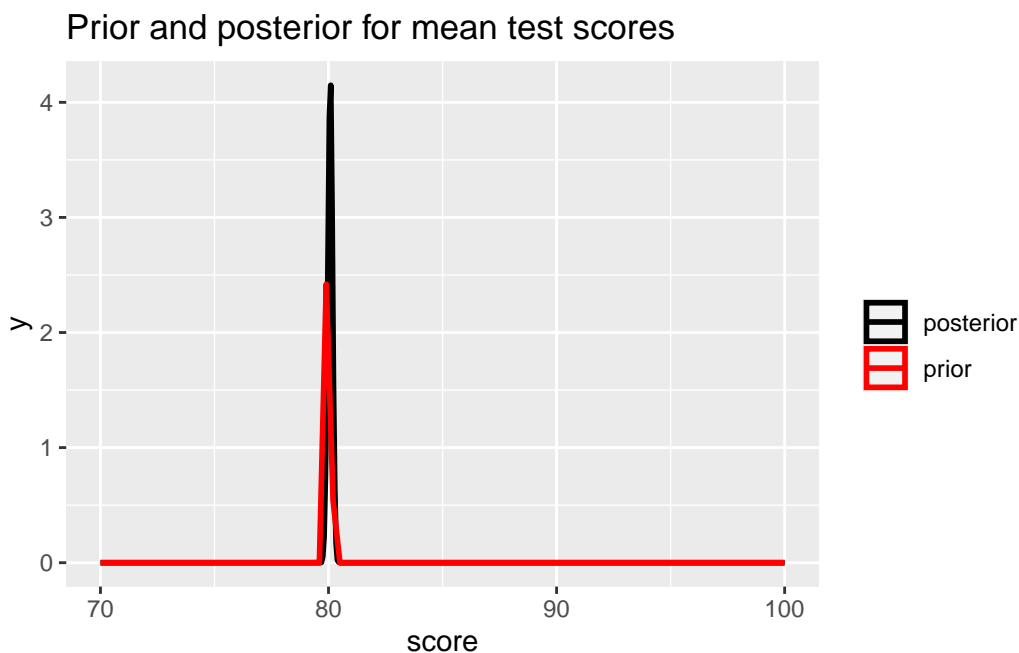
	mean	se_mean	sd	2.5%	25%	50%	75%	97.5%	n_eff
mu	80.07	0.00	0.10	79.87	80.00	80.07	80.13	80.26	583
sigma	21.43	0.03	0.76	20.02	20.93	21.42	21.90	23.07	911
lp__	-1548.41	0.06	1.11	-1551.33	-1548.84	-1548.03	-1547.65	-1547.39	322
Rhat									
mu	1.00								
sigma	1.00								
lp__	1.01								

Samples were drawn using NUTS(diag\_e) at Thu Feb 15 22:21:57 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).

```
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")
```



## 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.

```
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("labs/kids3.stan"),
             data = data,
             iter = 1000)
```



### Question 3

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

To show that the estimates of the intercept and slope are comparable for the results from `lm()`, we first fit a model using `lm()`. Then, we compare the fitted model with the `fit2` obtained from the bayes regression.

The the intercept term for linear regression is estimated as 77.548, while for bayes regression we have 77.96 as the mean intercept. The two values only differed by 0.412. Given the scale of the data, the means are comparable. We then look at the coefficient term for `mom_hs`. The coefficient for `lm()` is 11.771, for bayes regressoin, the coefficient is 11.25. The second one is just slightly greater than the one estimated from linear models. In addition, the estimated standard errors of these coefficients are similar as well. Thus, the results of the two models are comparable.

```
fit.lm<-lm(kid_score~mom_hs,data=kidiq)
summary(fit.lm)
```

Call:

```
lm(formula = kid_score ~ mom_hs, data = kidiq)
```

Residuals:

Min	1Q	Median	3Q	Max
-57.55	-13.32	2.68	14.68	58.45

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	77.548	2.059	37.670	< 2e-16 ***
mom_hs	11.771	2.322	5.069	5.96e-07 ***

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 19.85 on 432 degrees of freedom

Multiple R-squared: 0.05613, Adjusted R-squared: 0.05394

F-statistic: 25.69 on 1 and 432 DF, p-value: 5.957e-07

```
fit2
```

Inference for Stan model: anon\_model.

4 chains, each with iter=1000; warmup=500; thin=1;

post-warmup draws per chain=500, total post-warmup draws=2000.

	mean	se_mean	sd	2.5%	25%	50%	75%	97.5%
alpha	77.95	0.07	2.00	74.05	76.59	77.97	79.29	81.87
beta[1]	11.24	0.08	2.24	6.76	9.69	11.18	12.85	15.54
sigma	19.79	0.02	0.65	18.52	19.32	19.80	20.24	21.06
lp__	-1514.31	0.04	1.19	-1517.41	-1514.79	-1513.99	-1513.46	-1512.98

	n_eff	Rhat
alpha	855	1
beta[1]	830	1
sigma	1129	1
lp__	925	1

Samples were drawn using NUTS(diag\_e) at Thu Feb 15 22:22:49 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).

In addition, the variance estimated from the bayes regression is 20.27. The variance for `lm()` is 19.85, which is similar to the variance estimated from bayes regression as well.

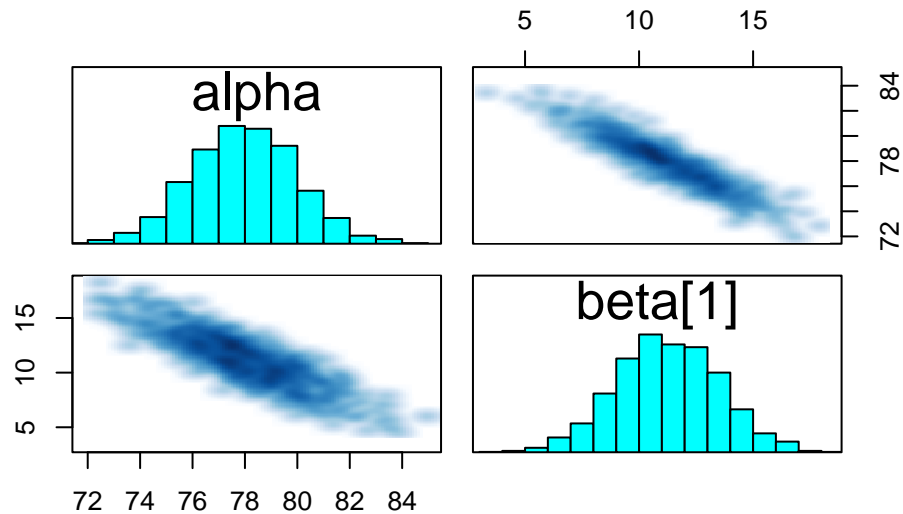
```
sigma2 <- sqrt(sum(fit.lm$residuals^2) / (length(y) - 2))
sigma2
```

```
[1] 19.85253
```

- 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?

From the `pairs` plot below, we observe a negative correlation between the intercept and the coefficient. This may not potentially be a problem, since in linear regression,  $\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$ , indicating the increase coefficient, the other one decreases.

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



#### 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.

```
kidiq$mom_iq_centered <- kidiq$mom_iq - mean(kidiq$mom_iq)

X <- as.matrix(kidiq[, c("mom_hs", "mom_iq_centered")])
K <- 2

data3 <- list(y = y,
              N = length(y),
              X = X,
              K = K
            )
fit3 <- stan(file = here("labs/kids3.stan"),
            data = data3,
            iter = 1000)
```

```
fit3
```

Inference for Stan model: anon\_model.

4 chains, each with iter=1000; warmup=500; thin=1;

post-warmup draws per chain=500, total post-warmup draws=2000.

mean	se_mean	sd	2.5%	25%	50%	75%	97.5%
------	---------	----	------	-----	-----	-----	-------

alpha	82.38	0.06	1.93	78.57	81.06	82.38	83.77	86.16
beta[1]	5.60	0.07	2.21	1.14	4.14	5.61	7.02	9.89
beta[2]	0.56	0.00	0.06	0.45	0.52	0.57	0.61	0.68
sigma	18.12	0.02	0.61	16.99	17.69	18.09	18.51	19.43
lp__	-1474.45	0.05	1.41	-1477.97	-1475.19	-1474.13	-1473.41	-1472.64
	n_eff	Rhat						
alpha	1053	1.01						
beta[1]	1032	1.01						
beta[2]	1445	1.00						
sigma	1513	1.00						
lp__	926	1.00						

Samples were drawn using NUTS(diag\_e) at Thu Feb 15 22:22:52 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).

The coefficient of the centered mom's IQ is estimated as 0.56. This value indicates that, if we keep other variables the same, on average, an unit increase in the mom's IQ increases the kids' score by 0.56.

## Question 5

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

From the output below, we observe that, the centered mom's IQ from `lm()` is estimated as 0.56 as well, with aligns with our bayes regression estimate.

```
fit3.lm <- lm(kid_score ~ mom_hs + mom_iq_centered, data=kidiq)
summary(fit3.lm)
```

Call:

```
lm(formula = kid_score ~ mom_hs + mom_iq_centered, data = kidiq)
```

Residuals:

Min	1Q	Median	3Q	Max
-52.873	-12.663	2.404	11.356	49.545

Coefficients:

Estimate	Std. Error	t value	Pr(> t )
----------	------------	---------	----------

```

(Intercept)      82.12214      1.94370  42.250 < 2e-16 ***
mom_hs           5.95012      2.21181   2.690 0.00742 **
mom_iq_centered  0.56391      0.06057   9.309 < 2e-16 ***
---

```

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 18.14 on 431 degrees of freedom

Multiple R-squared: 0.2141, Adjusted R-squared: 0.2105

F-statistic: 58.72 on 2 and 431 DF, p-value: < 2.2e-16

## Question 6

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

From density plot below, we observe that, for moms who have an IQ of 110, if the moms did not graduate from high school, their kids' IQ is approximately normally distributed with mean around 87, if the moms has graduated from highschool, the disribution of their kids IQ has a mean around 93.

```

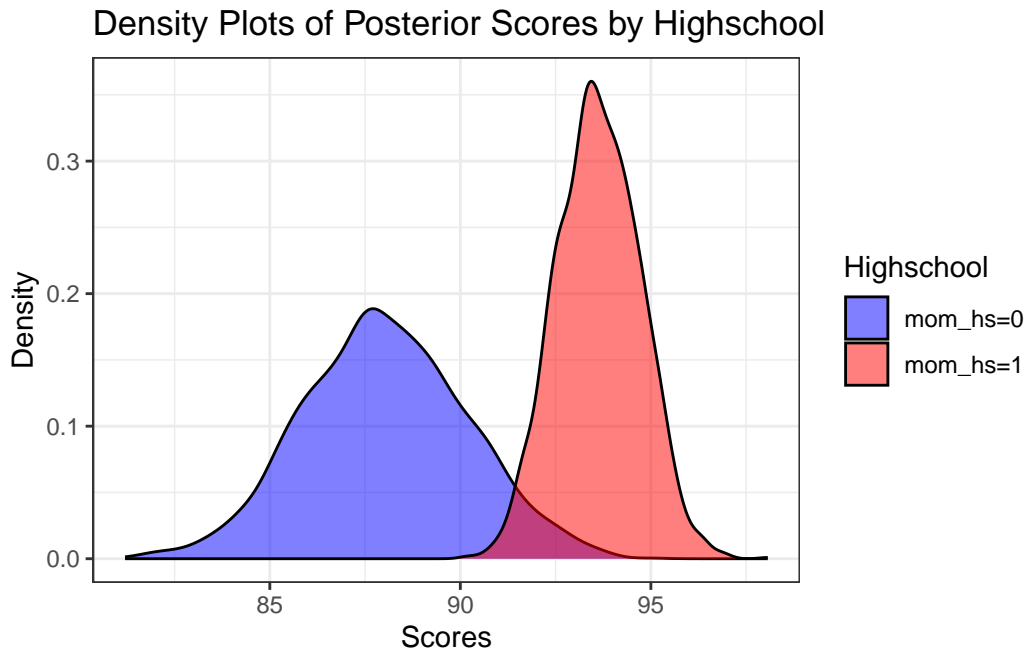
posterior_samples <- extract(fit3)
b0=posterior_samples$alpha
b1=posterior_samples$beta[,1]
b2=posterior_samples$beta[,2]
e=posterior_samples$sigma

posterior.hs0=b0+b2*(110-mean(kidiq$mom_iq))
posterior.hs1=b0+b1+b2*(110-mean(kidiq$mom_iq))

density<- data.frame(
  Scores = c(posterior.hs0, posterior.hs1),
  Highschool = rep(c("mom_hs=0", "mom_hs=1"), each = length(posterior.hs0)))

ggplot(density, aes(x = Scores, fill = Highschool)) +
  geom_density(alpha = 0.5) +
  labs(title = "Density Plots of Posterior Scores by Highschool", x = "Scores", y = "Density")

```



```
#a<-density|> filter(Highschool=="mom_hs=1")
#a$Scores |>mean()
```

### 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.

We observe that the estimated kids IQ are approximately normally distributed with a mean around 103 .

```
posterior.hs1.95=b0+b1+b2*(95-mean(kidiq$mom_iq))+e

ggplot()+geom_histogram(aes(x=posterior.hs1.95),fill="lightblue")+theme_bw()+labs(x="IQ",t
```

Distribution of kids' IQ whose mother graduated from high school

