

# HW 1 Solutions

Madeleine Livaudais

Fall 2025

## 7.2 Fake-data simulation and regression:

Simulate 100 data points from the linear model,  $y = a + bx + \text{error}$ , with  $a = 5$ ,  $b = 7$ , the values of  $x$  being sampled at random from a uniform distribution on the range  $[0, 50]$ , and errors that are normally distributed with mean 0 and standard deviation 3.

### 7.2a

Fit a regression line to these data and display the output.

```
a <- 5
b <- 7
sigma <- 3

x <- runif(n = 100, min = 0, max = 50)
n <- length(x)

y <- a + b*x + rnorm(n, 0, sigma)
sim_dat <- data.frame(x, y)

fit <- stan_glm(y ~ x, data = sim_dat)
```

```
##
## SAMPLING FOR MODEL 'continuous' NOW (CHAIN 1).
## Chain 1:
## Chain 1: Gradient evaluation took 0.001213 seconds
## Chain 1: 1000 transitions using 10 leapfrog steps per transition would take 12.13 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: 0.061 seconds (Warm-up)
## Chain 1: 0.044 seconds (Sampling)
## Chain 1: 0.105 seconds (Total)
## Chain 1:
##
## SAMPLING FOR MODEL 'continuous' NOW (CHAIN 2).
## Chain 2:
## Chain 2: Gradient evaluation took 8e-06 seconds
## Chain 2: 1000 transitions using 10 leapfrog steps per transition would take 0.08 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: 0.061 seconds (Warm-up)
## Chain 2: 0.041 seconds (Sampling)
## Chain 2: 0.102 seconds (Total)
## Chain 2:
##
## SAMPLING FOR MODEL 'continuous' 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: 0.051 seconds (Warm-up)

```

```

## Chain 3:          0.044 seconds (Sampling)
## Chain 3:          0.095 seconds (Total)
## Chain 3:
##
## SAMPLING FOR MODEL 'continuous' NOW (CHAIN 4).
## Chain 4:
## Chain 4: Gradient evaluation took 6e-06 seconds
## Chain 4: 1000 transitions using 10 leapfrog steps per transition would take 0.06 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: 0.052 seconds (Warm-up)
## Chain 4:          0.038 seconds (Sampling)
## Chain 4:          0.09 seconds (Total)
## Chain 4:

```

```
print(fit)
```

```

## stan_glm
## family:      gaussian [identity]
## formula:     y ~ x
## observations: 100
## predictors:  2
## -----
##               Median MAD_SD
## (Intercept)  4.9      0.6
## x            7.0      0.0
##
## Auxiliary parameter(s):
##               Median MAD_SD
## sigma  3.0      0.2
##
## -----
## * For help interpreting the printed output see ?print.stanreg
## * For info on the priors used see ?prior_summary.stanreg

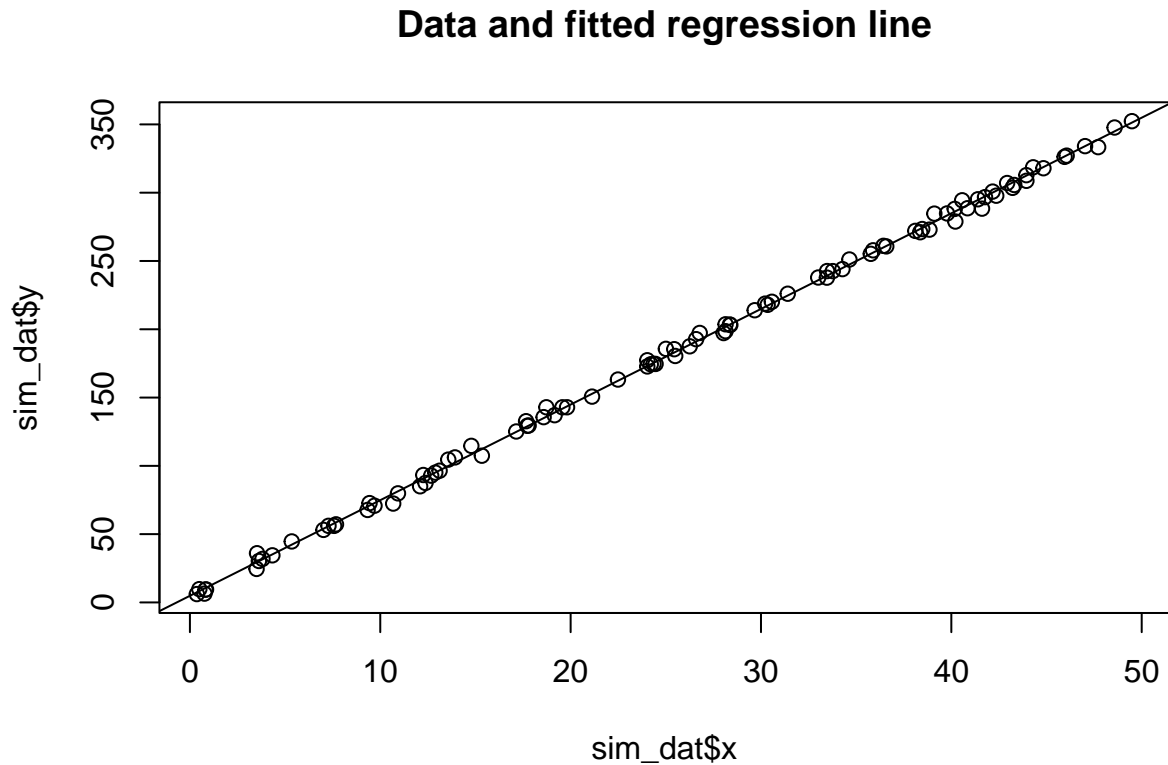
```

## 7.2b

Graph a scatterplot of the data and the regression line.

```
plot(sim_dat$x, sim_dat$y, main = "Data and fitted regression line")

a_hat <- coef(fit)[1]
b_hat <- coef(fit)[2]
abline(a_hat, b_hat)
```



#### 7.2c

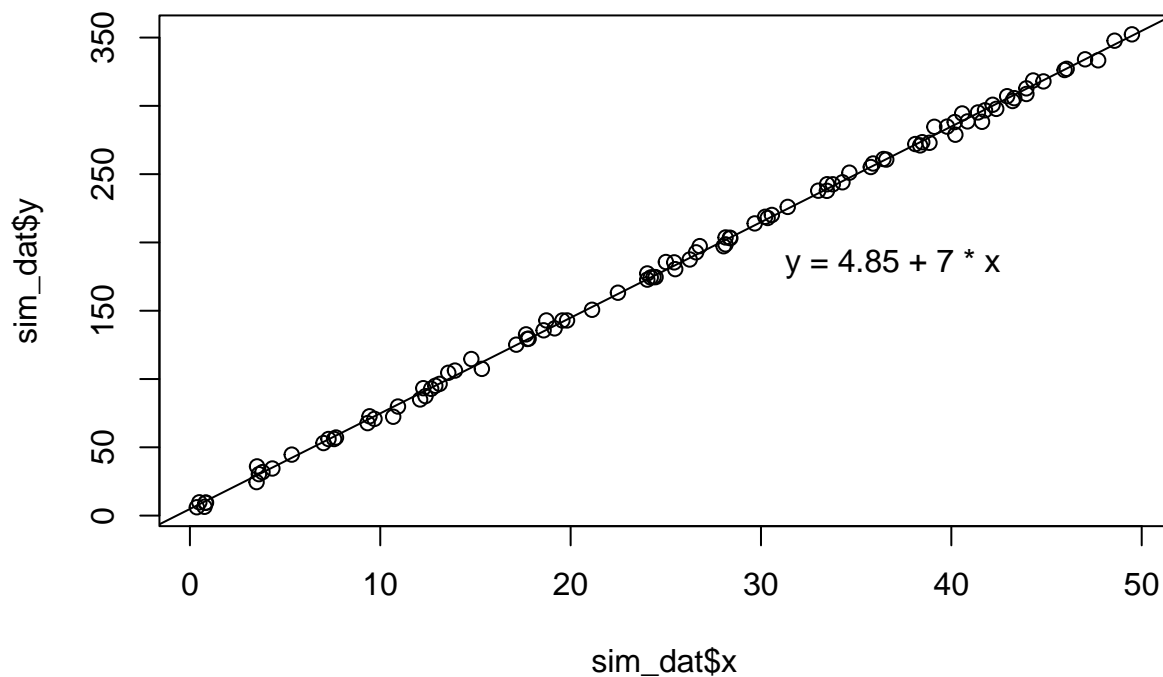
Use the `text` function in R to add the formula of the fitted line to the graph.

```
plot(sim_dat$x, sim_dat$y, main = "Data and fitted regression line")

a_hat <- coef(fit)[1]
b_hat <- coef(fit)[2]
abline(a_hat, b_hat)

x_bar <- mean(sim_dat$x)
text(x_bar, a_hat + b_hat*x_bar, paste("y =", round(a_hat, 2), "+", round(b_hat, 2), "* x"), adj = -.5)
```

## Data and fitted regression line



### 7.3 Fake-data simulation and fitting the wrong model:

Simulate 100 data points from the model  $y = a + bx + cx^2 + \text{error}$ , with the values of  $x$  being sampled at random from a uniform distribution on the range  $[0, 50]$ , errors that are normally distributed with mean 0 and standard deviation 3, and  $a, b, c$  chosen so that a scatterplot of the data shows a clear nonlinear curve.

#### 7.3 a

Fit a regression line `stan_glm(y ~ x)` to these data and display the output.

```
a <- 5
b <- 7
c <- 10
sigma <- 3

x <- runif(n = 100, min = 0, max = 50)
n <- length(x)

y <- a + b*x + c*x^2 + rnorm(n, 0, sigma)
sim_dat <- data.frame(x, y)

fit <- stan_glm(y ~ x, data = sim_dat)
```

```
##
```

```

## SAMPLING FOR MODEL 'continuous' NOW (CHAIN 1).
## Chain 1:
## Chain 1: Gradient evaluation took 1.6e-05 seconds
## Chain 1: 1000 transitions using 10 leapfrog steps per transition would take 0.16 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: 0.729 seconds (Warm-up)
## Chain 1:                0.031 seconds (Sampling)
## Chain 1:                0.76 seconds (Total)
## Chain 1:
##
## SAMPLING FOR MODEL 'continuous' NOW (CHAIN 2).
## Chain 2:
## Chain 2: Gradient evaluation took 8e-06 seconds
## Chain 2: 1000 transitions using 10 leapfrog steps per transition would take 0.08 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: 0.568 seconds (Warm-up)
## Chain 2:                0.032 seconds (Sampling)
## Chain 2:                0.6 seconds (Total)
## Chain 2:
##
## SAMPLING FOR MODEL 'continuous' 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 / 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: 0.522 seconds (Warm-up)
## Chain 3:                0.032 seconds (Sampling)
## Chain 3:                0.554 seconds (Total)
## Chain 3:
##
## SAMPLING FOR MODEL 'continuous' 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: 0.61 seconds (Warm-up)
## Chain 4:                0.032 seconds (Sampling)
## Chain 4:                0.642 seconds (Total)
## Chain 4:

```

```
print(fit)
```

```

## stan_glm
## family:      gaussian [identity]
## formula:     y ~ x
## observations: 100
## predictors:  2
## -----

```

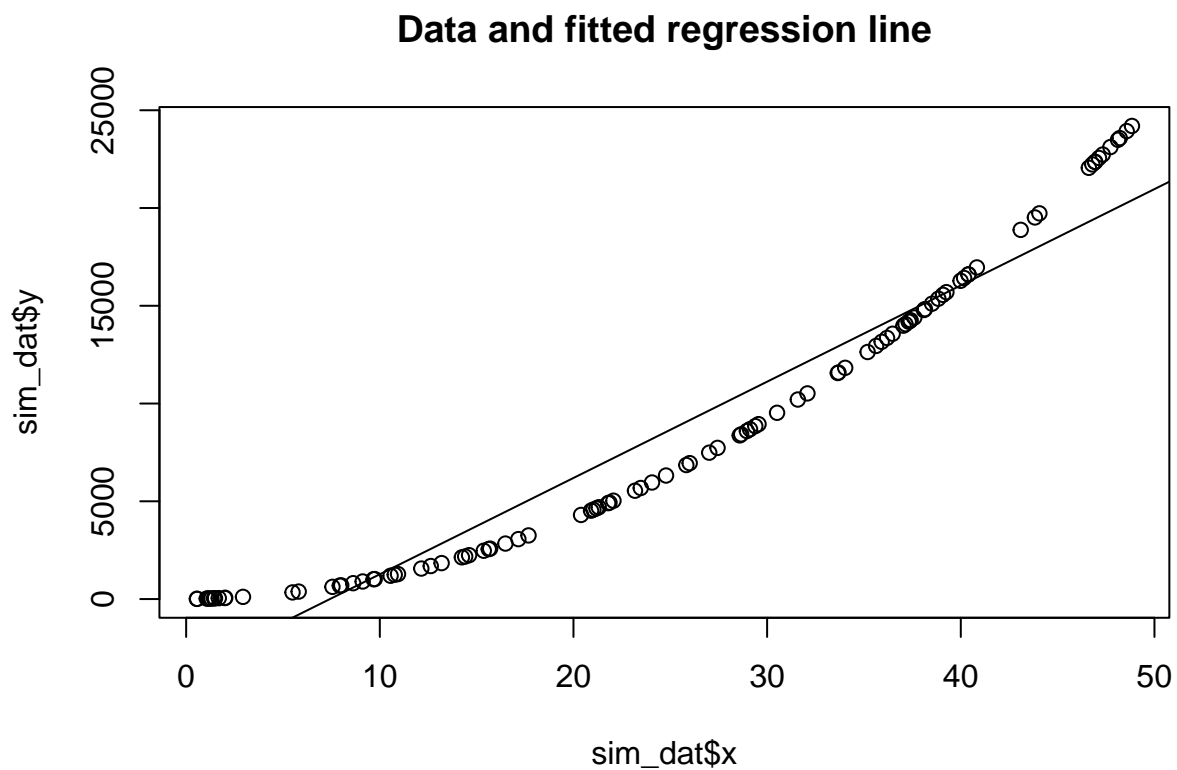
```
##           Median  MAD_SD
## (Intercept) -3656.0   385.9
## x           492.4    13.2
##
## Auxiliary parameter(s):
##           Median MAD_SD
## sigma 1974.3   144.5
##
## -----
## * For help interpreting the printed output see ?print.stanreg
## * For info on the priors used see ?prior_summary.stanreg
```

### 7.3b

Graph a scatterplot of the data and the regression line. This is the best-fit linear regression. What does “best-fit” mean in this context?

```
plot(sim_dat$x, sim_dat$y, main = "Data and fitted regression line")

a_hat <- coef(fit)[1]
b_hat <- coef(fit)[2]
abline(a_hat, b_hat)
```





## 7.6 Formulating comparisons as regression models:

Take the election forecasting model and simplify it by creating a binary predictor defined as  $x = 0$  if income growth is less than 2% and  $x = 1$  if income growth is more than 2%.

### 7.6a

Compute the difference in incumbent party's vote share on average, comparing those two groups of elections, and determine the standard error for this difference.

### 7.6b

Regress incumbent party's vote share on the binary predictor of income growth and check that the resulting estimate and standard error are the same as above.

## 8.8 Comparing `lm` and `stan_glm`:

Use simulated data to compare least squares estimation to default Bayesian regression:

### 8.8a

Simulate 100 data points from the model,  $y = 2 + 3x + \text{error}$ , with predictors  $x$  drawn from a uniform distribution from 0 to 20 and with independent errors drawn from the normal distribution with mean 0 and standard deviation 5. Fit the regression of  $y$  on  $x$  data using `lm` and `stan_glm` (using its default settings) and check that the two programs give nearly identical results.

```
a <- 2
b <- 3
sigma <- 5

x <- runif(n = 100, min = 0, max = 20)
n <- length(x)

y <- a + b*x + rnorm(n, 0, sigma)
sim_dat <- data.frame(x, y)

fit_1 <- stan_glm( y ~ x, data = sim_dat)
```

```
##
## SAMPLING FOR MODEL 'continuous' NOW (CHAIN 1).
## Chain 1:
## Chain 1: Gradient evaluation took 1.6e-05 seconds
## Chain 1: 1000 transitions using 10 leapfrog steps per transition would take 0.16 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: 0.032 seconds (Warm-up)
## Chain 1: 0.032 seconds (Sampling)
## Chain 1: 0.064 seconds (Total)
## Chain 1:
##
## SAMPLING FOR MODEL 'continuous' 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: 0.032 seconds (Warm-up)
## Chain 2: 0.031 seconds (Sampling)
## Chain 2: 0.063 seconds (Total)
## Chain 2:
##
## SAMPLING FOR MODEL 'continuous' 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:           0.033 seconds (Sampling)
## Chain 3:           0.066 seconds (Total)
## Chain 3:
##
## SAMPLING FOR MODEL 'continuous' NOW (CHAIN 4).
## Chain 4:
## Chain 4: Gradient evaluation took 8e-06 seconds
## Chain 4: 1000 transitions using 10 leapfrog steps per transition would take 0.08 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: 0.024 seconds (Warm-up)
## Chain 4:           0.03 seconds (Sampling)
## Chain 4:           0.054 seconds (Total)
## Chain 4:

```

```
print(fit_1)
```

```

## stan_glm
## family:      gaussian [identity]
## formula:     y ~ x
## observations: 100
## predictors:  2
## -----
##               Median MAD_SD
## (Intercept) 1.5      1.0
## x           3.1      0.1
##
## Auxiliary parameter(s):
##           Median MAD_SD
## sigma 5.2      0.4
##
## -----
## * For help interpreting the printed output see ?print.stanreg
## * For info on the priors used see ?prior_summary.stanreg

```

```
fit_2 <- lm( y ~ x, data = sim_dat)
print(fit_2)
```

```
##
## Call:
## lm(formula = y ~ x, data = sim_dat)
##
## Coefficients:
## (Intercept)          x
##      1.467      3.085
```

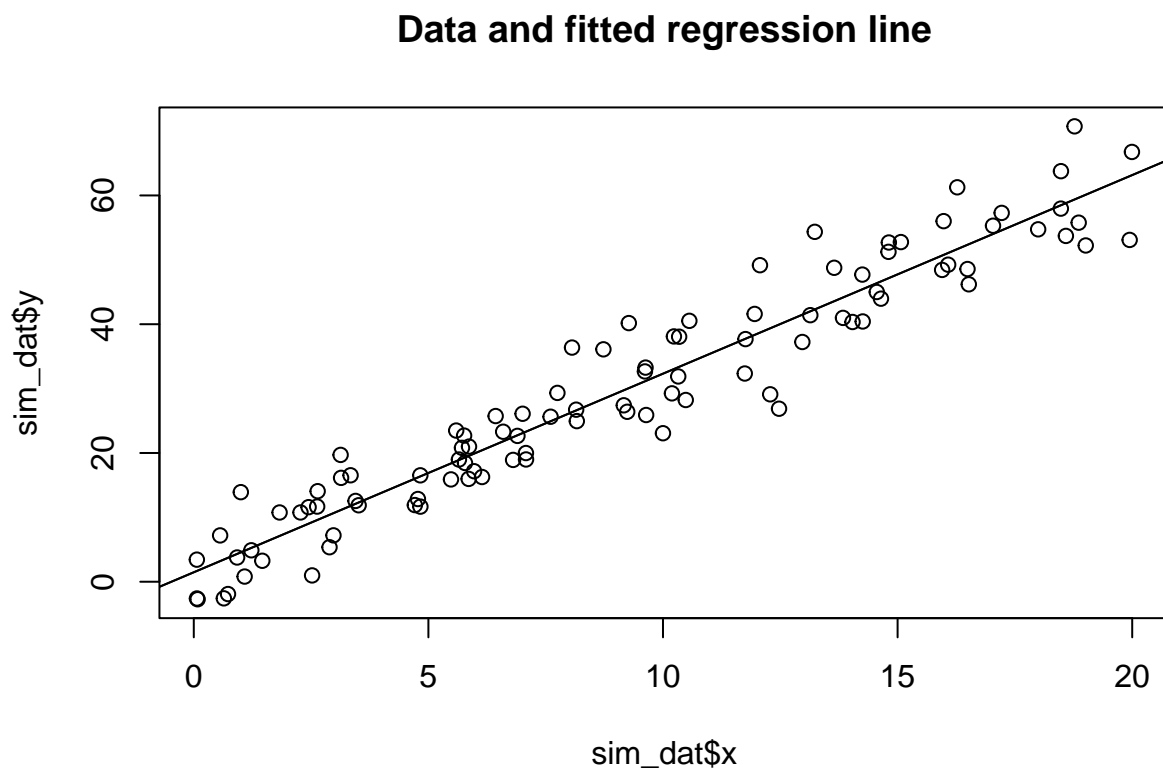
8.8b

Plot the simulated data and the two fitted regression lines.

```
plot(sim_dat$x, sim_dat$y, main = "Data and fitted regression line")

a_hat_1 <- coef(fit_1)[1]
b_hat_1 <- coef(fit_1)[2]
abline(a_hat_1, b_hat_1)

a_hat_2 <- coef(fit_2)[1]
b_hat_2 <- coef(fit_2)[2]
abline(a_hat_2, b_hat_2)
```



### 8.8c

Repeat the two steps above, but try to create conditions for your simulation so that `lm` and `stan_glm` give much different results.

## 10.1 Regression with interactions:

Simulate 100 data points from the model,  $y = b_0 + b_1x + b_2z + b_3xz + \text{error}$ , with a continuous predictor  $x$  and a binary predictor  $z$ , coefficients  $b = c(1, 2, -1, -2)$ , and errors drawn independently from a normal distribution with mean 0 and standard deviation 3, as follows. For each data point  $i$ , first draw  $z_i$ , equally likely to take on the values 0 and 1. Then draw  $x_i$  from a normal distribution with mean  $z_i$  and standard deviation 1. Then draw the error from its normal distribution and compute  $y_i$ .

### 10.1a

Display your simulated data as a graph of  $y$  vs  $x$ , using dots and circles for the points with  $z = 0$  and 1, respectively.

### 10.1b

Fit a regression predicting  $y$  from  $x$  and  $z$  with no interaction. Make a graph with the data and two parallel lines showing the fitted model.

### 10.1c

Fit a regression predicting  $y$  from  $x$ ,  $z$ , and their interaction. Make a graph with the data and two lines showing the fitted model.

## 10.2 Regression with interactions:

Here is the output from a fitted linear regression of outcome  $y$  on pre-treatment predictor  $x$ , treatment indicator  $z$ , and their interaction:

```
              Median MAD_SD
(Intercept)  1.2      0.2
x             1.6      0.4
z             2.7      0.3
x:z           0.7      0.5

Auxiliary parameter(s):
      Median MAD_SD
sigma  0.4      0.0
```

### 10.2a

Write the equation of the estimated regression line of  $y$  on  $x$  for the treatment group and the control group, and the equation of the estimated regression line of  $y$  on  $x$  for the control group.

## 10.2b

Graph with pen on paper the two regression lines, assuming the values of  $x$  fall in the range  $(0, 10)$ . On this graph also include a scatterplot of data (using open circles for treated units and dots for controls) that are consistent with the fitted model.

## 10.5 Regression modeling and prediction:

The folder `KidIQ` contains a subset of the children and mother data discussed earlier in the chapter. You have access to children's test scores at age 3, mother's education, and the mother's age at the time she gave birth for a sample of 400 children.

### 10.5a

Fit a regression of child test scores on mother's age, display the data and fitted model, check assumptions, and interpret the slope coefficient. Based on this analysis, when do you recommend mothers should give birth? What are you assuming in making this recommendation?

### 10.5b

Repeat this for a regression that further includes mother's education, interpreting both slope coefficients in this model. Have your conclusions about the timing of birth changed?

### 10.5c

Now create an indicator variable reflecting whether the mother has completed high school or not. Consider interactions between high school completion and mother's age. Also create a plot that shows the separate regression lines for each high school completion status group.

### 10.5d

Finally, fit a regression of child test scores on mother's age and education level for the first 200 children and use this model to predict test scores for the next 200. Graphically display comparisons of the predicted and actual scores for the final 200 children.

## 10.6 Regression models with interactions:

The folder `Beauty` contains data (use file `beauty.csv`) from Hamermesh and Parker (2005) on student evaluations of instructors' beauty and teaching quality for several courses at the University of Texas. The teaching evaluations were conducted at the end of the semester, and the beauty judgments were made later, by six students who had not attended the classes and were not aware of the course evaluations.

See also Felton, Mitchell, and Stinson (2003) for more on this topic.

### 10.6a

Run a regression using `beauty` (the variable `beauty`) to predict course evaluations (`eval`), adjusting for various other predictors. Graph the data and fitted model, and explain the meaning of each of the coefficients along with the residual standard deviation. Plot the residuals versus fitted values.

### 10.6b

Fit some other models, including beauty and also other predictors. Consider at least one model with interactions. For each model, explain the meaning of each of its estimated coefficients.

## 10.7 Predictive simulation for linear regression:

Take one of the models from the previous exercise.

### 10.7a

Instructor A is a 50-year-old woman who is a native English speaker and has a beauty score of -1. Instructor B is a 60-year-old man who is a native English speaker and has a beauty score of -0.5. Simulate 1000 random draws of the course evaluation rating of these two instructors. In your simulation, use `posterior_predict` to account for the uncertainty in the regression parameters as well as predictive uncertainty.

### 10.7b

Make a histogram of the difference between the course evaluations for A and B. What is the probability that A will have a higher evaluation?

## 10.8 How many simulation draws:

Take the model from Exercise 10.6 that predicts course evaluations from beauty and other predictors.

### 10.8a

Display and discuss the fitted model. Focus on the estimate and standard error for the coefficient of beauty.

### 10.8b

Compute the median and mad sd of the posterior simulations of the coefficient of beauty, and check that these are the same as the output from printing the fit.

### 10.8c

Fit again, this time setting `iter = 1000` in your `stan_glm` call. Do this a few times in order to get a sense of the simulation variability.

### 10.8d

Repeat the previous step, setting `iter = 100` and then `iter = 10`.

### 10.8e

How many simulations were needed to give a good approximation to the mean and standard error for the coefficient of beauty?