

Homework 3

Anna Cook

Disclaimer

A few things to keep in mind : 1) Use `set.seed()` to make sure that the document produces the same random simulation as when you ran the code. 2) Use `refresh=0` for any `stan_glm()` or stan-based model. `lm()` or non-stan models don't need this! 3) You can type outside of the r chunks and make new r chunks where it's convenient. Make sure it's clear which questions you're answering. 4) Even if you're not too confident, please try giving an answer to the text responses! 5) Please don't print data in the document unless the question asks. It's good for you to do it to look at the data, but not as good for someone trying to read the document later on. 6) Check your document before submitting! Please put your name where "name" is by the author!

4.1 Comparison of proportions

A randomized experiment is performed within a survey. 1000 people are contacted. Half the people contacted are promised a \$5 incentive to participate, and half are not promised an incentive. The result is a 50% response rate among the treated group and 40% response rate among the control group. Give an estimate and standard error of the average treatment effect.

```
# estimate for treatment effect:
0.5 - 0.4

## [1] 0.1

# standard error of the treatment effect:
sqrt((0.5 / sqrt(500))^2 + (0.49 / sqrt(500))^2)

## [1] 0.03130815
```

4.2 Choosing sample size

You are designing a survey to estimate the gender gap: the difference in support for a candidate among men and women. Assuming the respondents are a simple random sample of the voting population, how many people do you need to poll so that the standard error is less than 5 percentage points?

```
p1 <- 0.5
p2 <- 0.5
n <- ((p1 * (1-p1)) + (p2 * (1-p2))) / (0.05)^2

# for se of the difference less than 5%, n must be at least 200.
```

4.4 Designing an experiment

You want to gather data to determine which of two students is a better basketball shooter. You plan to have each student take N shots and then compare their shooting percentages. Roughly how large does N have to be for you to have a good chance of distinguishing a 30% shooter from a 40% shooter?

```
p1 <- 0.3
p2 <- 0.4
```

```
n <- ((p1 * (1-p1)) + (p2 * (1-p2))) / (0.05)^2

# for a standard error of the difference between shooters less than 0.05, n must be at least 180
```

4.6 Hypothesis testing

The following are the proportions of girl births in Vienna for each month in Girl births 1908 and 1909 (out of an average of 3900 births per month):

```
birthdata <- c(.4777,.4875,.4859,.4754,.4874,.4864,.4813,.4787,.4895,.4797,.4876,.4859,.4857,.4907,.5011)
```

The data are in the folder Girls. These proportions were used by von Mises (1957) to support a claim that that the sex ratios were less variable than would be expected under the binomial distribution. We think von Mises was mistaken in that he did not account for the possibility that this discrepancy could arise just by chance.

(a) Compute the standard deviation of these proportions and compare to the standard deviation that would be expected if the sexes of babies were independently decided with a constant probability over the 24-month period.

```
sd(birthdata)

## [1] 0.006409724

# sd = sqrt(p*(1-p)/n) = sqrt(0.25/24) = 0.1021
# the sd that would be expected with a constant probability is much higher than the sd of the observed
```

(b) The observed standard deviation of the 24 proportions will not be identical to its theoretical expectation. In this case, is this difference small enough to be explained by random variation? Under the randomness model, the actual variance should have a distribution with expected value equal to the theoretical variance, and proportional to a chi-square random variable with 23 degrees of freedom; see page 53.

```
0.1021 / 0.0064

## [1] 15.95312

qchisq(birthdata, 23)

## [1] 21.96609 22.12860 22.10202 21.92803 22.12694 22.11032 22.02571 21.98264
## [9] 22.16184 21.99920 22.13026 22.10202 22.09870 22.18180 22.35359 22.17515
## [17] 22.10368 22.18846 22.12195 21.88010 22.04063 22.12029 22.04229 22.29179

#the theoretical value is almost 16 times greater than the observed value, which is is higher than woul
```

5.5 Distribution of averages and differences

The heights of men in the United States are approximately normally distributed with mean 69.1 inches and standard deviation 2.9 inches. The heights of women are approximately normally distributed with mean 63.7 inches and standard deviation 2.7 inches. Let x be the average height of 100 randomly sampled men, and y be the average height of 100 randomly sampled women. In R, create 1000 simulations of $x - y$ and plot their histogram. Using the simulations, compute the mean and standard deviation of the distribution of $x - y$ and compare to their exact values.

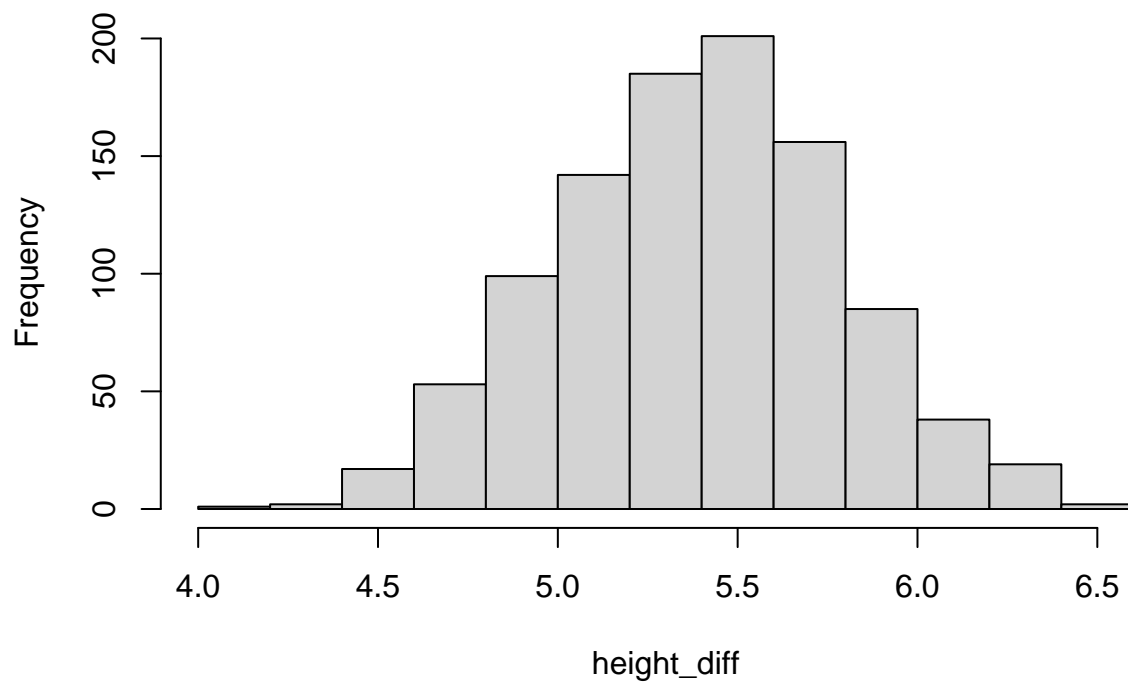
```
# simulation
set.seed(50)
```

```

n_sims <- 1000
height_diff <- rep(NA, n_sims)
for (s in 1:n_sims) {
  n <- 100
  men <- rnorm(n, 69.1, 2.9)
  women <- rnorm(n, 63.7, 2.7)
  x <- mean(men)
  y <- mean(women)
  height_diff[s] <- x-y
}
hist(height_diff)

```

Histogram of height_diff



```
mean(height_diff)
```

```
## [1] 5.38674
```

```
sd(height_diff)
```

```
## [1] 0.3928915
```

```
mad(height_diff)
```

```
## [1] 0.3962354
```

```
#actual values
```

```
mean_height_diff <- 69.1 - 63.7
```

```
mean_height_diff
```

```
## [1] 5.4
```

```
sd_height_diff <- sqrt((2.9)^2 + (2.7)^2)/sqrt(n)
```

```
sd_height_diff
```

```
## [1] 0.3962323
```

```
# From this output we can see that the simulated and exact values are very close to one another.
```

5.6 Propagation of uncertainty:

We use a highly idealized setting to illustrate the use of simulations in combining uncertainties. Suppose a company changes its technology for widget production, and a study estimates the cost savings at 5 dollars per unit, but with a standard error of 4 dollars. Furthermore, a forecast estimates the size of the market (that is, the number of widgets that will be sold) at 40 000, with a standard error of 10 000. Assuming these two sources of uncertainty are independent, use simulation to estimate the total amount of money saved by the new product (that is, savings per unit, multiplied by size of the market).

```
set.seed(30)
n_sims <- 1000
mean_savings <- rep(NA, n_sims)
for (s in 1:n_sims) {
  n <- 100
  savings_per_unit <- rnorm(n, 5, 4)
  market <- rnorm(n, 40000, 10000)
  total_savings <- savings_per_unit * market
  mean_savings[s] <- mean(total_savings)
}
mean(mean_savings)
```

```
## [1] 200540.7
```

5.8 Coverage of confidence intervals:

On page 15 there is a discussion of an experimental study of an education-related intervention in Jamaica, in which the point estimate of the treatment effect, on the log scale, was 0.35 with a standard error of 0.17. Suppose the true effect is 0.10—this seems more realistic than the point estimate of 0.35—so that the treatment on average would increase earnings by 0.10 on the log scale. Use simulation to study the statistical properties of this experiment, assuming the standard error is 0.17.

(a) Simulate 1000 independent replications of the experiment assuming that the point estimate is normally distributed with mean 0.10 and standard deviation 0.17.

```
set.seed(100)
n_sims <- 1000
mean_p_hat <- rep(NA, n_sims)
sd_p_hat <- rep(NA, n_sims)
for (s in 1:n_sims) {
  n <- 127
  p_hat <- rnorm(n, 0.10, 0.17)
  mean_p_hat[s] <- mean(p_hat)
  sd_p_hat[s] <- sd(p_hat)
}
```

(b) For each replication, compute the 95% confidence interval. Check how many of these intervals include the true parameter value.

```
n <- 127
se <- sd_p_hat/sqrt(n)
lower <- mean_p_hat + qt(0.025, n-1)*se
```

```
upper <- mean_p_hat + qt(0.975, n-1)*se
int_95 <- data.frame(lower, upper)
sum(int_95$lower <= 0.10 & int_95$upper >= 0.10)
```

```
## [1] 944
```

(c) Compute the average and standard deviation of the 1000 point estimates; these represent the mean and standard deviation of the sampling distribution of the estimated treatment effect.

```
mean(mean_p_hat)
```

```
## [1] 0.1003613
```

```
mean(sd_p_hat)
```

```
## [1] 0.1692568
```

5.9 Coverage of confidence intervals after selection on statistical significance:

Take your 1000 simulations from Exercise 5.8, and select just the ones where the estimate is statistically significantly different from zero. Compute the average and standard deviation of the selected point estimates. Compare these to the result from Exercise 5.8.

```
sum(int_95$lower > 0)
```

```
## [1] 1000
```

```
# Because none of the confidence intervals contain 0, this means that all 1000 estimates are statistical
```

9.8 Simulation for decision analysis:

An experiment is performed to measure the efficacy of a television advertising program. The result is an estimate that each minute spent on a national advertising program will increase sales by 500,000 dollars, and this estimate has a standard error of 200,000 dollars. Assume the uncertainty in the treatment effect can be approximated by a normal distribution. Suppose ads cost 300,000 dollars per minute. What is the expected net gain for purchasing 20 minutes of ads? What is the probability that the net gain is negative?

```
set.seed(100)
sales <- rnorm(100, 500000, 200000)
cost <- 300000
min <- 20
net_gain <- (sales*min) - (cost*min)
mean(net_gain)
```

```
## [1] 4011650
```

```
z <- (0 - mean(sales) / sd(sales))
pnorm(z)
```

```
## [1] 0.007100705
```

10.3 Checking statistical significance:

In this exercise and the next, you will simulate two variables that are statistically independent of each other to see what happens when we run a regression to predict one from the other. Generate 1000 data points from a normal distribution with mean 0 and standard deviation 1 by typing `var1 <- rnorm(1000,0,1)` in R.

Generate another variable in the same way (call it var2). Run a regression of one variable on the other. Is the slope coefficient “statistically significant”? We do not recommend summarizing regressions in this way, but it can be useful to understand how this works, given that others will do so.

```
set.seed(80)
var1 <- rnorm(1000,0,1)
var2 <- rnorm(1000,0,1)
df <- data.frame(var1,var2)
fit <- lm(var1 ~ var2, data = df)
summary(fit)
```

```
##
## Call:
## lm(formula = var1 ~ var2, data = df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3.1675 -0.5933  0.0316  0.6198  3.3329
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -0.03463    0.03112  -1.113   0.266
## var2         0.03345    0.03057   1.094   0.274
##
## Residual standard error: 0.9838 on 998 degrees of freedom
## Multiple R-squared:  0.001198,    Adjusted R-squared:  0.0001969
## F-statistic: 1.197 on 1 and 998 DF,  p-value: 0.2742
```

#no, the slope coefficient is not significant. The p-value = 0.274.

10.4 Simulation study of statistical significance:

Continuing the previous exercise, run a simulation repeating this process 100 times. This can be done using a loop. From each simulation, save the z-score (the estimated coefficient of var1 divided by its standard error). If the absolute value of the z-score exceeds 2, the estimate is “statistically significant.” To perform this computation, we start by creating an empty vector of z-scores filled with missing values (NAs). Another approach is to start with `z_scores <- numeric(length=100)`, which would set up a vector of zeroes. In general, however, we prefer to initialize with NAs, because then when there is a bug in the code, it sometimes shows up as NAs in the final results, alerting us to the problem.

How many of these 100 z-scores exceed 2 in absolute value, thus achieving the conventional level of statistical significance?

Here is code to perform the simulation:

This chunk will have `eval=FALSE`. If you want it to run, please copy it to a new chunk, or remove `eval=FALSE`!

```
z_scores <- rep(NA,100)
for(k in 1:100) {
  var1 <- rnorm(1000,0,1)
  var2 <- rnorm(1000,0,1)
  fake <- data.frame(var1,var2)
  fit <- stan_glm(var2 ~ var1,data=fake,refresh=0)
  z_scores[k] <- coef(fit)[2] / se(fit)[2]
}

sum(z_scores > abs(2))
```

```
## [1] 3
```

```
# 3 z-scores exceed an absolute value of 2
```

11.3 Coverage of confidence intervals:

Consider the following procedure:

- Set $n = 100$ and draw n continuous values x_i uniformly distributed between 0 and 10. Then simulate data from the model $y_i = a + bx_i + \text{error}_i$, for $i = 1, \dots, n$, with $a = 2$, $b = 3$, and independent errors from a normal distribution.
- Regress y on x . Look at the median and mad sd of b . Check to see if the interval formed by the median ± 2 mad sd includes the true value, $b = 3$.
- Repeat the above two steps 1000 times.

```
set.seed(60)
n <- 100
a <- 2
b <- 3
x <- runif(n, 0, 10)
y <- a + b*x + rnorm(n,0,1)
df <- data.frame(x,y)
fit <- stan_glm(y ~ x, data = df, refresh = 0)
```

(a) True or false: the interval should contain the true value approximately 950 times. Explain your answer.

#####true. With normal distribution, 95% of the data fall within ± 2 mad sd from the median (or sd from the mean). So if we calculate a 95% confidence interval for 1000 simulations, we should see that the true value is included in the interval 95% of the time, or 950 times.

(b) Same as above, except the error distribution is bimodal, not normal. True or false: the interval should contain the true value approximately 950 times. Explain your answer.

#####In a bimodal distribution, the standard deviation would be larger, so the confidence intervals would be wider, so we would expect more than 950 of the confidence intervals to contain the true value.

Optional:

11.6 Fitting a wrong model:

Suppose you have 100 data points that arose from the following model: $y = 3 + 0.1 x_1 + 0.5 x_2 + \text{error}$, with independent errors drawn from a t distribution with mean 0, scale 5, and 4 degrees of freedom. We shall explore the implications of fitting a standard linear regression to these data.

#####(a) Simulate data from this model. For simplicity, suppose the values of x_1 are simply the integers from 1 to 100, and that the values of x_2 are random and equally likely to be 0 or 1. In R, you can define `x_1 <- 1:100`, simulate `x_2` using `rbinom`, then create the linear predictor, and finally simulate the random errors in `y` using the `rt` function. Fit a linear regression (with normal errors) to these data and see if the 68% confidence intervals for the regression coefficients (for each, the estimates ± 1 standard error) cover the true values.

(b) Put the above step in a loop and repeat 1000 times. Calculate the confidence coverage for the 68% intervals for each of the three coefficients in the model.

11.9 Leave-one-out cross validation:

Use LOO to compare different models fit to the beauty and teaching evaluations example from Exercise 10.6:

###(a) Discuss the LOO results for the different models and what this implies, or should imply, for model choice in this example.

(b) Compare predictive errors pointwise. Are there some data points that have high predictive errors for all the fitted models?

11.10 K-fold cross validation:

Repeat part (a) of the previous example, but using 5-fold cross validation:

###(a) Randomly partition the data into five parts using the sample function in R.

(b) For each part, re-fitting the model excluding that part, then use each fitted model to predict the outcomes for the left-out part, and compute the sum of squared errors for the prediction.

(c) For each model, add up the sum of squared errors for the five steps in (b). Compare the different models based on this fit.

(d) Not in the textbook, but if you're curious, compare your results to `kfold()` or `cv.glm()`!