

Problem Set 10

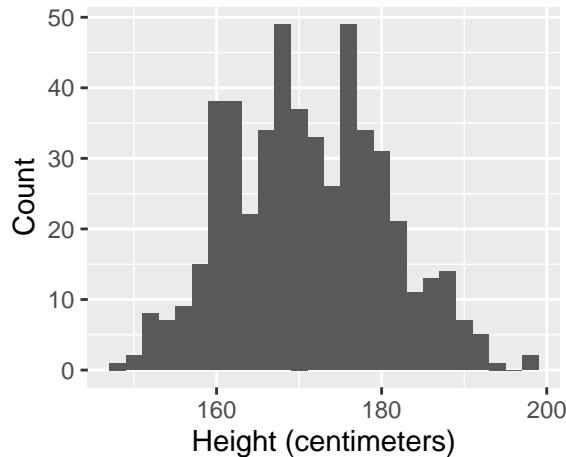
Inference for Means

- The following data set¹ was privatized through a combination of cell suppression and generalization. What level of k -anonymity does it demonstrate with respect to the variables Age, Zip Code, and Nationality? (i.e. what is k ?)

Age	Zip Code	Nationality	Disease
<30	130 **	*	Heart Disease
<30	130 **	*	Heart Disease
<30	130 **	*	Viral Infection
<30	130 **	*	Viral Infection
≥40	1485 *	*	Cancer
≥40	1485 *	*	Heart Disease
≥40	1485 *	*	Viral Infection
≥40	1485 *	*	Viral Infection
3 *	130 **	*	Cancer
3 *	130 **	*	Cancer
3 *	130 **	*	Cancer
3 *	130 **	*	Cancer

- Heights of adults.** Researchers studying anthropometry collected body measurements, as well as age, weight, height and gender, for 507 physically active individuals. Summary statistics for the distribution of heights (measured in centimeters), along with a histogram, are provided below.² (Heinz et al. 2003)

Min	Q1	Median	Mean	Q3	Max	SD	IQR
147.2	163.8	170.3	171.1	177.8	198.1	9.4	14



- What is the point estimate for the average height of active individuals? What about the median?
- What is the point estimate for the standard deviation of the heights of active individuals? What about the IQR?

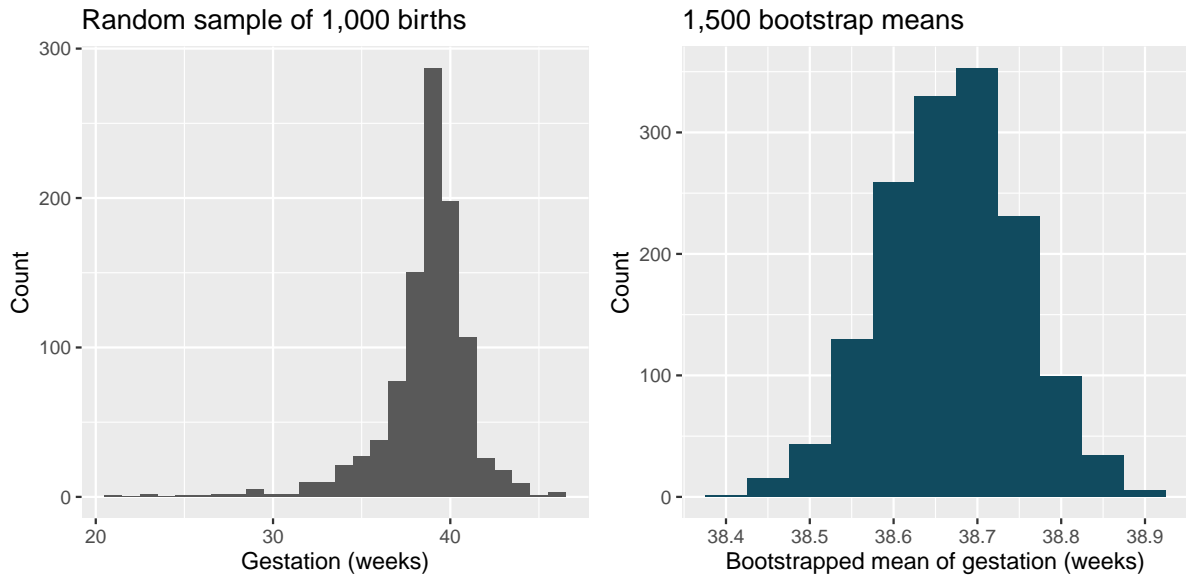
¹Widodo et al. (2019). *Privacy Preserving Data Publishing with Multiple Sensitive Attributes based on Overlapped Slicing*. Information. 10(12):362.

²The `bdims` data used in this exercise can be found in the `openintro` R package.

- c. Is a person who is 1m 80cm (180 cm) tall considered unusually tall? And is a person who is 1m 55cm (155cm) considered unusually short? Explain your reasoning.
 - d. The researchers take another random sample of physically active individuals. Would you expect the mean and the standard deviation of this new sample to be the ones given above? Explain your reasoning.
 - e. The sample means obtained are point estimates for the mean height of all active individuals, if the sample of individuals is equivalent to a simple random sample. What measure do we use to quantify the variability of such an estimate? Compute this quantity using the data from the original sample under the condition that the data are a simple random sample.
3. ***Heights of adults vs. kindergartners.** Heights of 507 physically active individuals have a mean of 171 centimeters and a standard deviation of 9.4 centimeters.³ (Heinz et al. 2003)
- a. Would the standard deviation of the heights of a few hundred kindergartners be bigger or smaller than 9.4cm? Explain your reasoning.
 - b. Suppose many samples of size 100 adults is taken and, separately, many samples of size 100 kindergartners are taken. For each of the many samples, the average height is computed. Which set of sample averages would have a larger standard error of the mean, the adult sample averages or the kindergartner sample averages?
4. **Possible bootstrap samples.** Consider a simple random sample of the following observations: 47, 4, 92, 47, 12, 8. Which of the following could be a possible bootstrap samples from the observed data above? If the set of values could not be a bootstrap sample, indicate why not.
- a. 47, 47, 47, 47, 47, 47
 - b. 92, 4, 13, 8, 47, 4
 - c. 92, 47, 12
 - d. 8, 47, 12, 12, 8, 4, 92
 - e. 12, 4, 8, 8, 92, 12
5. **Length of gestation, confidence interval.** Every year, the United States Department of Health and Human Services releases to the public a large dataset containing information on births recorded in the country. This dataset has been of interest to medical researchers who are studying the relation between habits and practices of expectant mothers and the birth of their children. In this exercise we work with a random sample of 1,000 cases from the dataset released in 2014. The length of pregnancy, measured in weeks, is commonly referred to as gestation. The histograms below show the distribution of lengths of gestation from the random sample of 1,000 births (on the left) and the distribution of bootstrapped means of gestation from 1,500 different bootstrap samples (on the right).⁴

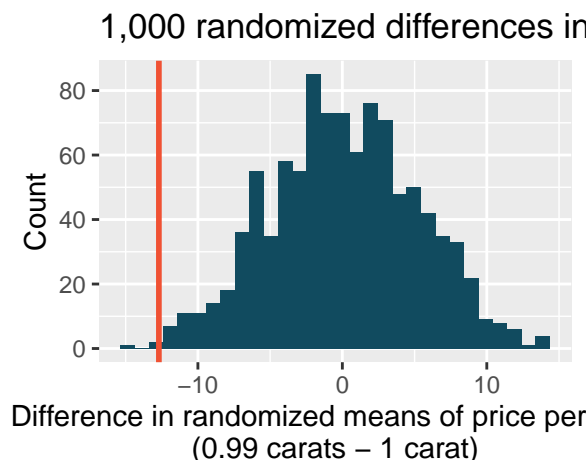
³The `bdims` data used in this exercise can be found in the `openintro` R package.

⁴The `births14` data used in this exercise can be found in the `openintro` R package.



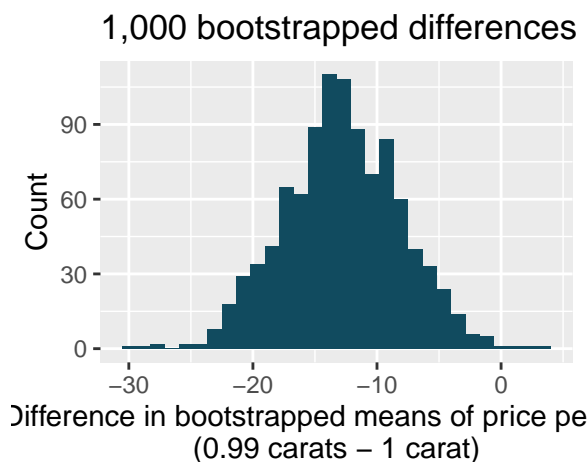
- Given the bootstrap sampling distribution for the sample mean, find an approximate value for the standard error of the mean.
 - By looking at the bootstrap sampling distribution (1,500 bootstrap samples were taken), find an approximate 99% bootstrap percentile confidence interval for the true average gestation length in the population from which the data were randomly sampled. Provide the interval as well as a one-sentence interpretation of the interval.
6. **Sleep habits of New Yorkers.** New York is known as “the city that never sleeps.” A random sample of 25 New Yorkers were asked how much sleep they get per night. Statistical summaries of these data are shown below. In this exercise you’ll use two mathematical approximation techniques to build a confidence interval.
- | n | Mean | SD | Min | Max |
|----|------|------|------|------|
| 25 | 7.73 | 0.77 | 6.17 | 9.78 |
- Use the t distribution to construct a 95% confidence interval for the population mean of hours sleep per night for all New Yorkers (see *CI for Numerical Data* slides or Ch. 19 in the textbook).
 - What are you assuming about the population distribution of hrs sleep per night when using the t distribution? Does that seem reasonable in this setting?
 - If you’re uncomfortable with the assumption from (b), you could alternatively rely upon the Central Limit Limit. Build a second confidence interval for the mean, this time using the normal distribution.
 - Under what conditions will the normal distribution provide a good approximation to the sampling distribution of the sample mean?
 - Based on the two confidence intervals that you’ve made, is 8 hours / night a plausible value for the population mean for all New Yorkers?
7. ***Diamonds, randomization test.** The prices of diamonds go up as the carat weight increases, but the increase is not smooth. For example, the difference between the size of a 0.99 carat diamond and a 1 carat diamond is undetectable to the naked human eye, but the price of a 1 carat diamond tends to be much higher than the price of a 0.99 diamond. In this question we use two random samples of diamonds, 0.99 carats and 1 carat, each sample of size 23, and randomize the carat weight to the price values in order compare the average prices of the diamonds to a null distribution. In order to be able to

compare equivalent units, we first divide the price for each diamond by 100 times its weight in carats. That is, for a 0.99 carat diamond, we divide the price by 99. or a 1 carat diamond, we divide the price by 100. The randomization distribution (with 1,000 repetitions) below describes the null distribution of the difference in sample means (of price per carat) if there really was no difference in the population from which these diamonds came.⁵ (Wickham 2016)



Using the randomization distribution of the difference in average price per carat (1,000 randomizations were run), conduct a hypothesis test to evaluate if there is a difference between the prices per carat of diamonds that weigh 0.99 carats and diamonds that weigh 1 carat. Make sure to state your hypotheses clearly and interpret your results in context of the data. (Wickham 2016)

8. ***Diamonds, bootstrap interval.** We have data on two random samples of diamonds: one with diamonds that weigh 0.99 carats and one with diamonds that weigh 1 carat. Each sample has 23 diamonds. Provided below is a histogram of bootstrap differences in means of price per carat of diamonds that weight 0.99 carats and diamonds that weigh 1 carat. (Wickham 2016)



Using the bootstrap distribution, create a (rough) 95% bootstrap percentile confidence interval for the true population difference in prices per carat of diamonds that weigh 0.99 carats and 1 carat. Interpret the interval in the context of the this problem.

Heinz, G., L. J. Peterson, R. W. Johnson, and C. J. Kerk. 2003. "Exploring Relationships in Body Dimensions." *Journal of Statistics Education* 11 (2). http://www.openintro.org/redirect.php?go=textbook-body_dim_2003.

⁵The `diamonds` data used in this exercise can be found in the `ggplot2` R package.

Wickham, Hadley. 2016. *ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York. <https://ggplot2.tidyverse.org>.