

FCT NOVA

PROJECT REPORT

# Bootstrap and Jackknife

*Leonard Storcks*

Course by Dr. Regina BISPO

January 12, 2024

## **Abstract**

This project is all about the usage of resampling in statistics, e.g. to estimate standard errors of estimators. We will employ bootstrap techniques for estimation and testing in different contexts as well as the Jackknife.

# Contents

<b>1</b>	<b>Bootstrap Methods</b>	<b>1</b>
1.1	Caveats of Bootstrap . . . . .	1
1.2	Application of Bootstrap for Estimating the Difference between Means (task 1a) . . . . .	3
1.3	Bootstrap Confidence Intervals . . . . .	6
1.3.1	Standard Normal Bootstrap Confidence Interval . . . . .	6
1.3.2	Percentile Bootstrap Confidence Interval . . . . .	7
1.3.3	Basic Bootstrap Confidence Interval . . . . .	7
1.3.4	Bootstrap Confidence Intervals for Task 1b . . . . .	7
1.4	Bootstrap Hypothesis Testing . . . . .	9
1.4.1	Bootstrap Hypothesis Test on the Difference of Means I: Theory . . .	11
1.4.2	Bootstrap Hypothesis Test on the Difference of Means I: Example Application . . . . .	12
1.4.3	Bootstrap Hypothesis Test on the Difference of Means III: Task 1c . .	13
1.5	Bootstrapping on Linear Regression Coefficients (task 2) . . . . .	16
1.5.1	Linear Fit on the Geysir Data (task 2a) . . . . .	16
1.6	Bootstrap estimates by resampling cases (task 2b) . . . . .	17
1.7	Standard normal bootstrap confidence intervals (task 2c) . . . . .	19
1.8	Testing on $H_0 : \beta_0 = 0$ or rather $H_0 : \beta_0 = c$ (task 2d) . . . . .	19
<b>2</b>	<b>Jackknife and Further Bootstrap Applications</b>	<b>22</b>
2.1	Bias correction and standard error estimation using Bootstrap and Jackknife (task 3a, 3b) . . . . .	22
2.2	Non-parametric bootstrap method to test the hypothesis $H_0 : X \sim \mathcal{N}(2, 1)$ (task 3c) . . . . .	26
	<b>References</b>	<b>29</b>

# 1 Bootstrap Methods

At the heart of bootstrap lies random sampling with replacement to better assess statistical estimates (assign measures of accuracy<sup>1</sup>) - we use a sample as an estimate of a population and bootstrap samples from the sample as estimates of the sampling distribution. Consider we have observed a sample  $X_{\text{sample}} = (x_1, \dots, x_n)$  and are interested in a population parameter  $\theta$  which we can estimate as  $\hat{\theta}$  from a sample. The non-parametric bootstrap method is given by

1.  $\forall b \in \{1, \dots, B\}$  (where  $B$  is the number of bootstrap samples): Generate a bootstrap sample  $X^{(b)} = (x_1^{(b)}, \dots, x_n^{(b)})$  by sampling from  $X$  with replacement and calculate  $\hat{\theta}^{(b)}$  on the bootstrap sample.
2. This yields an empirical distribution of estimates  $\hat{\theta}^{(1)}, \dots, \hat{\theta}^{(B)}$  from which we can calculate the mean  $\hat{\theta}^*$  and standard error  $\text{SE}(\hat{\theta}^*)$ .

with the bootstrap estimate and the estimated standard error given by

$$\begin{aligned}\hat{\theta}^* &= \frac{1}{B} \sum_{b=1}^B \hat{\theta}^{(b)} \\ \hat{\text{SE}}(\hat{\theta}) &= \text{SE}(\hat{\theta}^*) = \sqrt{\frac{1}{B-1} \sum_{b=1}^B (\hat{\theta}^{(b)} - \hat{\theta}^*)^2}\end{aligned}\tag{1}$$

The bootstrap estimate of the bias  $\hat{\text{bias}}(\hat{\theta}) = E[\hat{\theta}] - \theta$  of the estimator  $\hat{\theta}$  is given by

$$\hat{\text{bias}}(\hat{\theta}) = \hat{\theta}^* - \hat{\theta}\tag{2}$$

## 1.1 Caveats of Bootstrap

Beware, that bootstrap is not magic - it will not retrieve information not present in a sample or fix a bad sample. Also, the bootstrap sample is generally centered at the observed statistic, not the population parameter (e.g.  $\bar{X}$  not  $\mu$ ) - using bootstrap we will not improve on  $\bar{X}$  (Hesterberg, 2015, section 2.3 (also see there for exceptions)). And if there is bias (as estimated above), the bootstrap estimate is indeed not a bias-corrected estimate - if  $\hat{\theta}^*$  is greater than  $\hat{\theta}$ , the bias corrected estimate  $\bar{\theta} = \hat{\theta} - \hat{\text{bias}}$  should be less than  $\hat{\theta}$ , and this

---

<sup>1</sup>For the mean e.g. a measure of accuracy is the standard error for which we have the explicit formula  $\text{SE} = \sqrt{\frac{s^2}{n}}$ ,  $\bar{X}_{\text{sample}} = \frac{1}{n} \sum_{i=1}^n x_i$ ,  $s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{X}_{\text{sample}})^2$ ,  $X_{\text{sample}} = (x_1, \dots, x_n)$ . For most other statistics (e.g. the median) we do not have such an explicit formula. See Efron and Tibshirani, 1994, chapter 1.

kind of bias correction also has pitfalls<sup>2</sup> (Efron and Tibshirani, 1994, chapter 10.6, p. 138). Simply put, if from resampling a sample we get a positive bias then to get to the unbiased one we have to go "two steps back",  $\hat{\theta} = \hat{\theta}^* - 2 \cdot (\hat{\theta}^* - \hat{\theta}) = 2\hat{\theta} - \hat{\theta}^*$  (eq. 10.41 in Efron and Tibshirani, 1994).

In figure 1 we illustrate that as the sample carries distributional information so do the bootstrap samples. In figure 2 we illustrate the difference between the sample estimate, bootstrap estimate and corrected estimate for a biased estimator, the standard deviation. Here the bias comes from the fact that in the calculation of the standard deviation we also estimate the mean from which the points in the sample naturally deviate less than from the true mean, so the standard deviation is (with the biased estimator) underestimated.

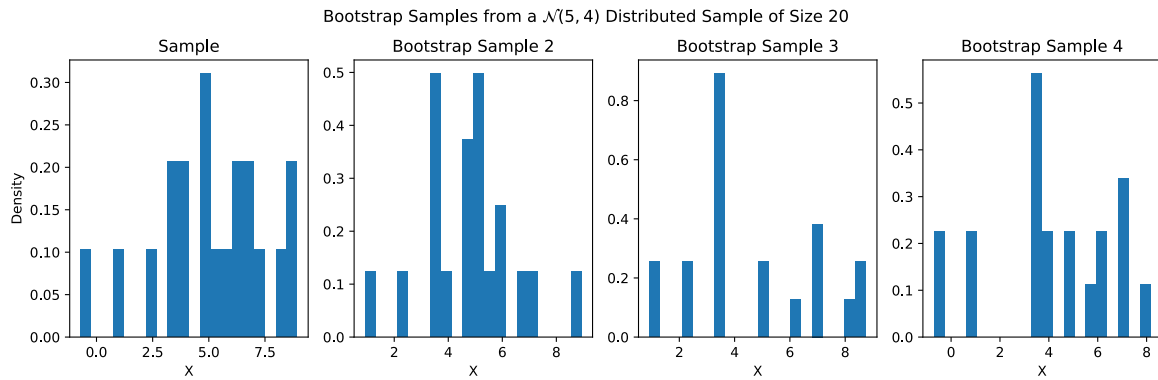


Figure 1: Bootstrap Samples from a sample.

---

<sup>2</sup> $\bar{\theta}$  might have a much larger standard error than  $\hat{\theta}$ .

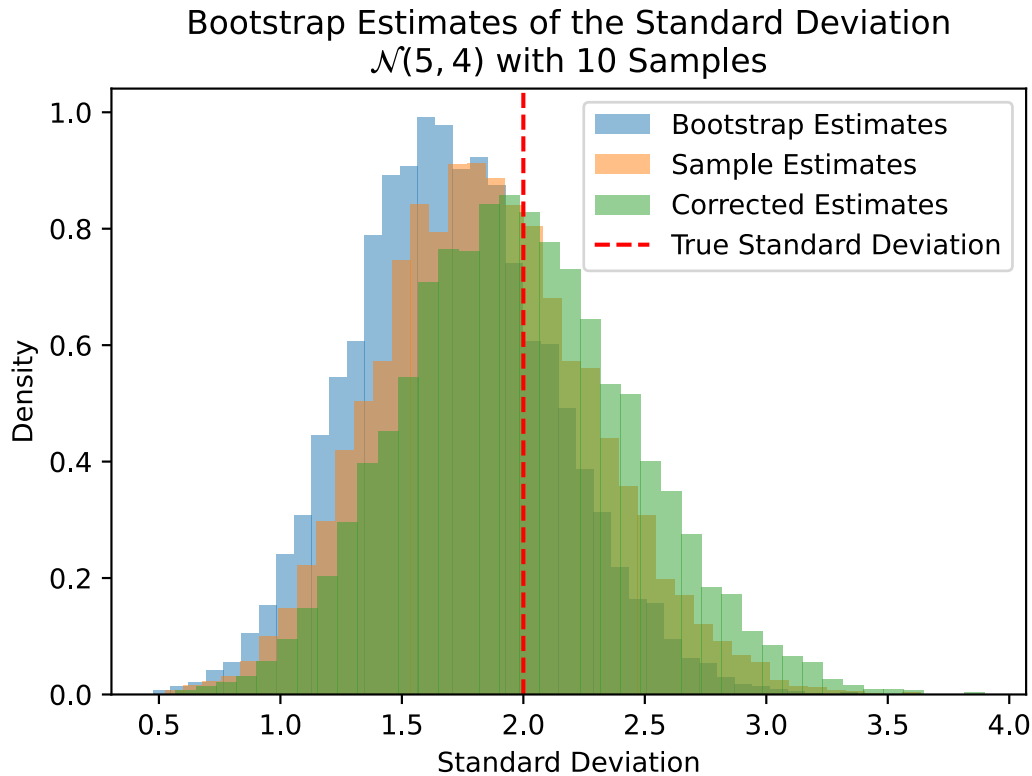


Figure 2: Illustration of the difference between the sample estimate, bootstrap estimate and corrected estimate for the biased estimate of the standard deviation.

## 1.2 Application of Bootstrap for Estimating the Difference between Means (task 1a)

In the following consider survival times of patients with breast cancer (sample 1) and stomach cancer (sample 2).

$$\begin{aligned} D^{(1)} &= (107, 353, 1764, 667, 990, 78, 667, 44, 9, 27), \quad n_1 = 10 \\ D^{(2)} &= (374, 253, 812, 246, 95, 367, 251, 309, 594, 826, 593, 97), \quad n_2 = 12 \end{aligned} \tag{3}$$

We will work with the logged lifetimes  $X^{(1)}, X^{(2)}, x_j^{(k)} = \log d_j^{(k)}$  (see e.g. Overduin, 2004) and present kernel density estimates of both the original and log-transformed data in figure 3. Note that in the sample 1 we have an outlier at 1764 which strongly pulls the mean of the non-logged data to higher values. The log-transforms crunches down especially the outlier, shifting the mean to lower values.

As  $\bar{X}^{(1)}$  is smaller than  $\bar{X}^{(2)}$  (see figure 3), we will consider  $\mu_2 - \mu_1$  instead of  $\mu_1 - \mu_2$  throughout the task contrary to the task description. This seems to be more in alignment of what the later test for  $H_0 : \mu_2 - \mu_1 \leq 0.25$  tries to get at: While  $\bar{X}^{(2)} - \bar{X}^{(1)} > 0.25$  we

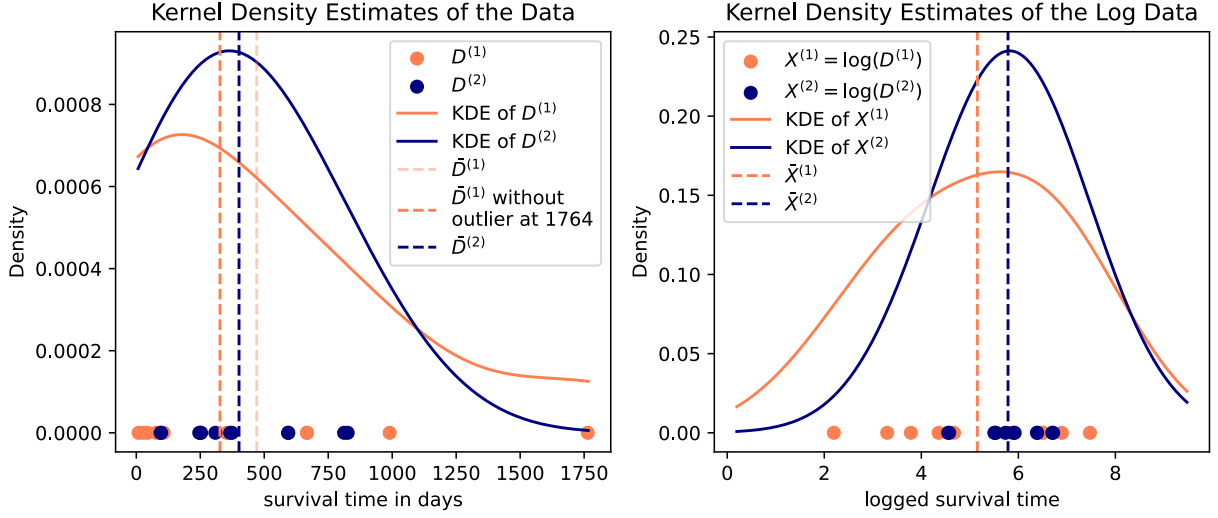


Figure 3: Kernel Density Estimates of the original and log-transformed data.

cannot reject  $H_0$  fundamentally as of the large stray in the distributions (see figure 3 again).

Let us come to the task and give a bootstrap-corrected estimate of  $\mu_2 - \mu_1$  with corresponding standard error.

The bootstrap procedure here is

1.  $\forall b \in \{1, \dots, B\}$ : Generate bootstrap samples  $X^{(1,b)}$  and  $X^{(2,b)}$  by sampling respectively  $n_1$  and  $n_2$  values from  $X^{(1)}$  and  $X^{(2)}$  with replacement and calculate  $t^{(b)} = \bar{X}^{(2,b)} - \bar{X}^{(1,b)}$  on the bootstrap samples.
2. This yields an empirical distribution of estimates  $t^{(1)}, \dots, t^{(B)}$  from which we can calculate the mean  $\hat{t}^*$  and standard error  $\text{SE}(t^*)$ .

The estimate of the bias is given by  $\hat{\text{bias}}(\hat{t}) = \hat{t}^* - \hat{t}$  with  $\hat{t} = \bar{X}^{(2)} - \bar{X}^{(1)}$ , so we get the bias corrected estimate  $\bar{t} = \hat{t} - \hat{\text{bias}}(\hat{t}) = 2\hat{t} - \hat{t}^*$ .

Note that by the bias correction we introduce a new source of variance (the variance of the bias) which we do not account for in the standard error. Using further bootstrapping we can estimate the standard error of  $\bar{t}$ , but if  $\hat{\text{bias}}(\hat{t})$  is small compared to  $\hat{\text{SE}}(\hat{t})$  it is safer to use  $\hat{t}$  than  $\bar{t}$  (see Efron and Tibshirani, 1994, chapter 10.6, p. 138).

For the case at hand the estimates are given in table 1 and the code to obtain them in code-snippet 1. The bootstrap distribution is shown in figure 4. As of the small bias compared to the standard error we use  $\hat{t}$  as our estimate of  $\mu_2 - \mu_1$  with standard error  $\hat{\text{SE}}(\hat{t})$ .

```

1 bootstrap_mean_difference <- function(sample1, sample2, B = 1000) {
2   n1 <- length(sample1)
3   n2 <- length(sample2)
4   # our statistic is a difference of means
5   t <- function(a, b) {
6     mean(b) - mean(a)
7   }
8   # point estimate
9   t_hat <- t(sample1, sample2)
10  # bootstrap samples
11  t_bs <- replicate(B, t(sample(sample1, n1, replace = TRUE),
    ↪ sample(sample2, n2, replace = TRUE)))
12  # bootstrap estimate
13  t_hat_star <- mean(t_bs)
14  # standard error
15  se <- sd(t_bs)
16  # bias
17  bias <- t_hat_star - t_hat
18  # bias-corrected estimate
19  t_hat_star_bc <- 2 * t_hat - t_hat_star
20  return(list("t_hat" = t_hat, "t_bs" = t_bs, "t_hat_star" =
    ↪ t_hat_star, "se" = se, "bias" = bias, "t_hat_star_bc" =
    ↪ t_hat_star_bc))
21 }
22
23 # Define the data
24 sample1 <- c(107, 353, 1764, 667, 990, 78, 667, 44, 9, 27)
25 sample2 <- c(374, 253, 812, 246, 95, 367, 251, 309, 594, 826, 593, 97)
26
27 # convert to log scale
28 sample1 <- log(sample1)
29 sample2 <- log(sample2)
30
31 res <- bootstrap_mean_difference(sample1, sample2)

```

Code-Snippet 1: Bootstrap estimates of  $\mu_2 - \mu_1$  for the samples  $X^{(1)}$  and  $X^{(2)}$  for  $B = 1000$ .



point estimate $\hat{t}$	$\hat{\text{SE}}(\hat{t})$	bootstrap estimate $\hat{t}^*$	$\hat{\text{bias}}(\hat{t})$	bias-corrected estimate $\bar{t}$
$\approx 0.6318$	$\approx 0.574$	$\approx 0.6355$	$\approx 0.0037$	$\approx 0.6281$

Table 1: Estimates regarding  $\mu_2 - \mu_1$  for the samples  $X^{(1)}$  and  $X^{(2)}$  for  $B = 1000$ .

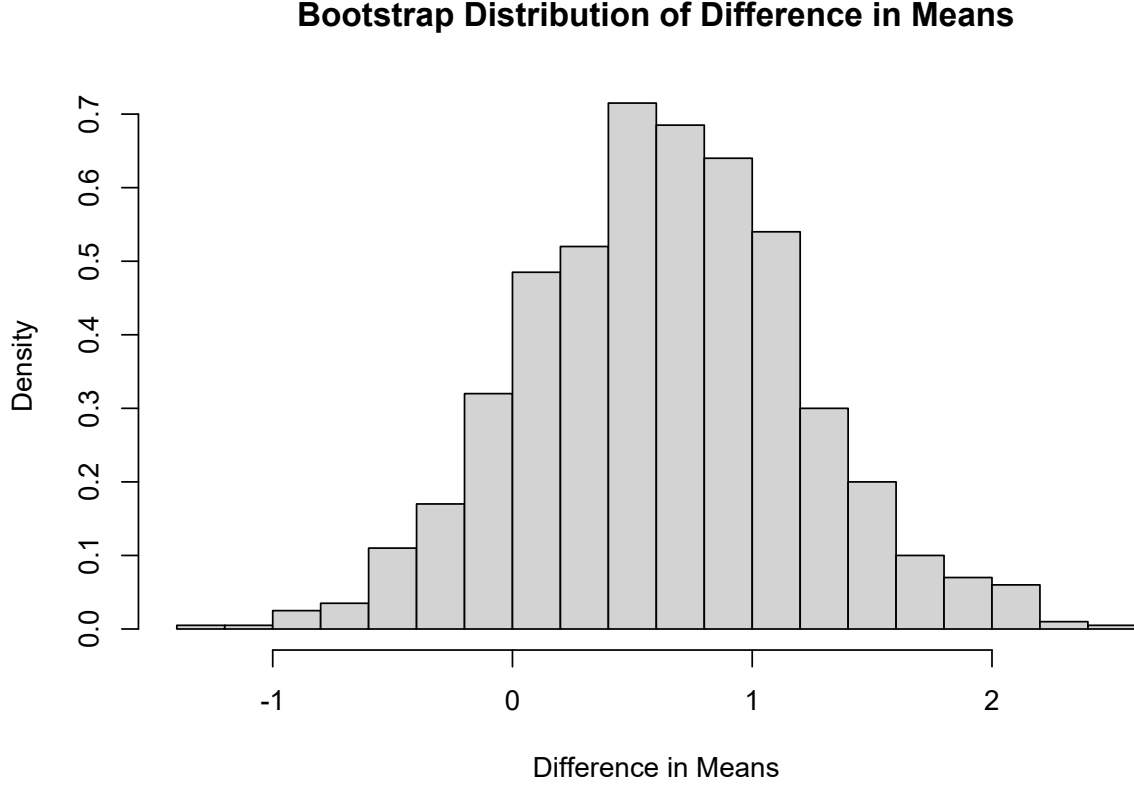


Figure 4: Bootstrap distribution of the difference in means of the bootstrap samples.

### 1.3 Bootstrap Confidence Intervals

Note in general that bootstrap distributions tend to be narrow on average, so the bootstrap confidence intervals often under-cover (Hesterberg, 2015).

#### 1.3.1 Standard Normal Bootstrap Confidence Interval

If  $\hat{\theta}$  is a sample mean, then for large sample sizes, we can apply the central limit theorem

$$\frac{\hat{\theta} - \theta}{\hat{\text{SE}}(\hat{\theta})} \underset{a}{\sim} N(0, 1) \quad (4)$$

so based on the bootstrap estimate of the standard error  $\hat{\text{SE}}(\hat{\theta})$  we can calculate a  $1 - \alpha$ -confidence interval for  $\theta$  as

$$\hat{\theta} \pm z_{1-\alpha/2} \hat{\text{SE}}(\hat{\theta}), \quad z_{1-\alpha/2} \text{ is the } 1 - \alpha/2 \text{ quantile of } N(0, 1) \quad (5)$$

assuming that  $\hat{\theta}$  is approximately normally distributed,  $\hat{\theta}$  is unbiased and  $\hat{\text{SE}}(\hat{\theta})$  is a good estimate of  $\text{SE}(\hat{\theta})$ .

### 1.3.2 Percentile Bootstrap Confidence Interval

Here, we use quantiles based on the empirical distribution of the bootstrap samples, so a  $1 - \alpha$ -confidence interval for  $\theta$  is given by

$$\begin{aligned} & \left[ \hat{\theta}_{\alpha/2}^*, \hat{\theta}_{1-\alpha/2}^* \right], \quad \hat{\theta}_{\alpha/2}^* \text{ is the } \alpha/2 \text{ quantile of } \hat{\theta}^*, \\ & \text{with } \hat{\theta}_{1-\alpha/2}^*, \hat{\theta}_{\alpha/2}^* \text{ being the quantiles of the empirical distribution } \left\{ \hat{\theta}^{(b)} \right\}_{b=1, \dots, B} \end{aligned} \quad (6)$$

### 1.3.3 Basic Bootstrap Confidence Interval

Akin to the bias-corrected estimate  $\bar{\theta}$ , we get a basic bootstrap confidence interval for  $\theta$  as

$$\begin{aligned} & \left[ 2\hat{\theta} - \hat{\theta}_{1-\alpha/2}^*, 2\hat{\theta} - \hat{\theta}_{\alpha/2}^* \right], \quad \hat{\theta}_{\alpha/2}^* \text{ is the } \alpha/2 \text{ quantile of } \hat{\theta}^*, \\ & \text{with } \hat{\theta}_{1-\alpha/2}^*, \hat{\theta}_{\alpha/2}^* \text{ being the quantiles of the empirical distribution } \left\{ \hat{\theta}^{(b)} \right\}_{b=1, \dots, B} \end{aligned} \quad (7)$$

### 1.3.4 Bootstrap Confidence Intervals for Task 1b

The  $\alpha = 0.05$ -confidence intervals for  $\mu_2 - \mu_1$  are given in table 2 and the code to obtain them in code-snippet 2 and 3. As expected based on the large stray in the distributions (see figure 3 again), the confidence intervals for the difference in the means are large, so our estimate of  $\mu_2 - \mu_1$  might not reflect the true value very well (even the sign might be wrong).

standard normal	percentile	basic
$[-0.49, 1.76]$	$[-0.49, 1.87]$	$[-0.61, 1.76]$

Table 2: Bootstrap confidence intervals for  $\mu_2 - \mu_1$  for the samples  $X^{(1)}$  and  $X^{(2)}$  for  $B = 1000$ .

```
1  standard_normal_interval <- function(SE_hat, t_hat, alpha = 0.05) {  
2    # SE_hat: standard error estimate  
3    # t_hat: point estimate  
4    # alpha: significance level  
5    # returns: confidence interval  
6  
7    # critical value  
8    z <- qnorm(1 - alpha / 2)  
9    # lower bound  
10   lower <- t_hat - z * SE_hat  
11   # upper bound  
12   upper <- t_hat + z * SE_hat  
13   return(c(lower, upper))  
14 }  
15  
16 percentile_boots_interval <- function(t_bs, alpha = 0.05) {  
17   # t_bs: bootstrap samples  
18   # alpha: significance level  
19   # returns: confidence interval  
20  
21   # critical values  
22   lower <- quantile(t_bs, alpha / 2)  
23   upper <- quantile(t_bs, 1 - alpha / 2)  
24   return(c(lower, upper))  
25 }  
26  
27 basic_boots_interval <- function(t_bs, t_hat, alpha = 0.05) {  
28   # t_bs: bootstrap samples  
29   # t_hat: point estimate  
30   # alpha: significance level  
31   # returns: confidence interval  
32  
33   # critical values  
34   lower <- 2 * t_hat - quantile(t_bs, 1 - alpha / 2)  
35   upper <- 2 * t_hat - quantile(t_bs, alpha / 2)  
36   return(c(lower, upper))  
37 }
```

Code-Snippet 2: General functions for calculating bootstrap confidence intervals.

```

1  # Define the data
2  sample1 <- c(107, 353, 1764, 667, 990, 78, 667, 44, 9, 27)
3  sample2 <- c(374, 253, 812, 246, 95, 367, 251, 309, 594, 826, 593,
4             ↪ 97)
5
6  # convert to log scale
7  sample1 <- log(sample1)
8  sample2 <- log(sample2)
9
10 res <- bootstrap_mean_difference(sample1, sample2)
11
12 # print and calculate confidence intervals
13 print(paste("standard normal interval:",
14             ↪ standard_normal_interval(res$se, res$t_hat)))
15 print(paste("percentile boots interval:",
16             ↪ percentile_boots_interval(res$t_bs)))
17 print(paste("basic boots interval:", basic_boots_interval(res$t_bs,
18             ↪ res$t_hat)))

```

Code-Snippet 3: Bootstrap confidence intervals for  $\mu_2 - \mu_1$  for the samples  $X^{(1)}$  and  $X^{(2)}$  for  $B = 1000$ .

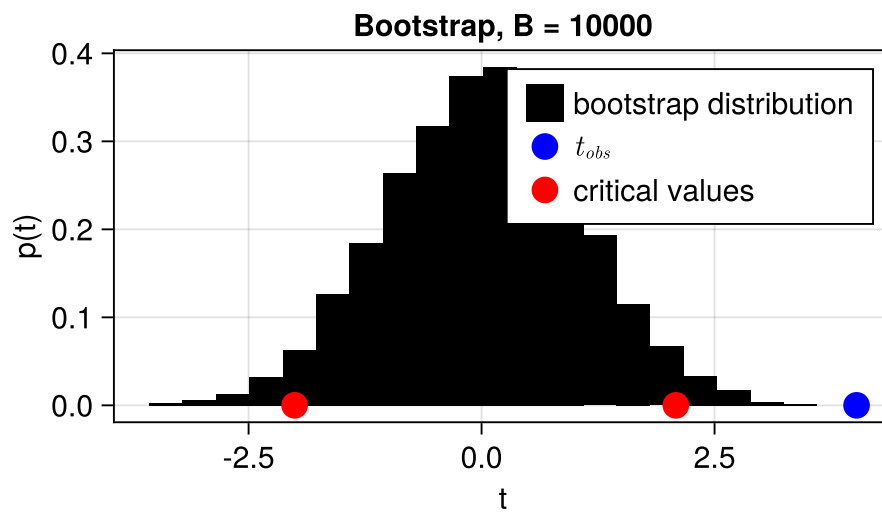
## 1.4 Bootstrap Hypothesis Testing

Consider we want to test a null hypothesis  $H_0$  and we have a test statistic  $T$  where observed values of  $T$  under  $H_0$  can be calculated for a sample as  $t_{\text{obs}}(X_{\text{sample}})$ . In a setting where we make assumptions on the distributions of the test statistic we would test based on where our observed value falls in the distribution of the test statistic under  $H_0$ .

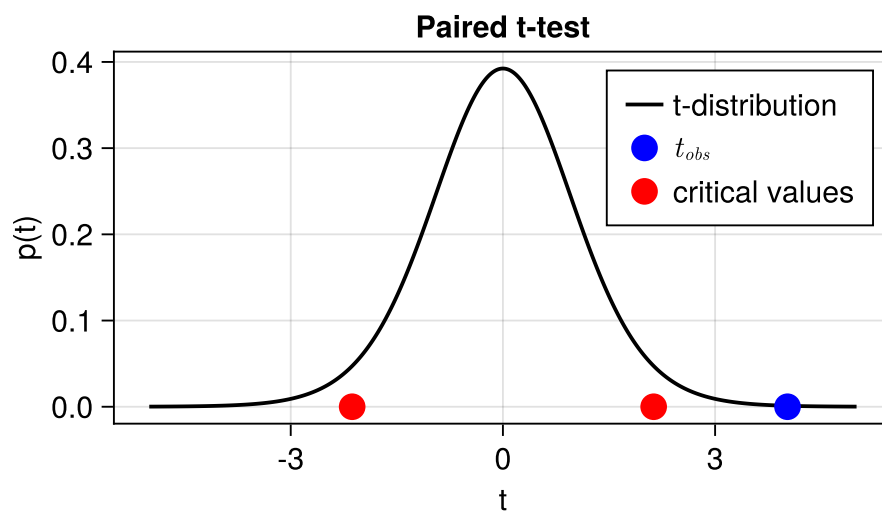
In the bootstrap setting we can test based on an empirical distribution based on resampling (see figure 5a), adapted to follow  $H_0$  if this is not the case by default. Consider for instance for a sample  $X_{\text{sample}} = (x_1, \dots, x_n)$  we want to test the null hypothesis  $H_0 : \mu = \mu_0$  vs  $H_1 : \mu \neq \mu_0$  (one-sample test). We then base our bootstrap samples on  $z_i = x_i - \bar{X}_{\text{sample}} + \mu_0$  and use the test statistic  $t^{(b)} = t^*(Z^{(b)}) = \frac{\bar{Z}^{(b)} - \mu_0}{s_Z^{(b)}/\sqrt{n}}$  where  $Z^{(b)} = (z_1^{(b)}, \dots, z_n^{(b)})$  and  $s_Z^{(b)} =$

$$\sqrt{\frac{1}{n-1} \sum_{i=1}^n \left( z_i^{(b)} - \bar{Z}^{(b)} \right)^2} \text{ and } t_{\text{obs}} = t^*(X_{\text{sample}}).$$

Here we assume that as the mean varies, the distributions are just translated versions of each other (translation family) where if our  $x_i$  are lifetimes it might make sense to use logged lifetimes as they are more likely to satisfy a translation or normal family assumption (Efron and Tibshirani, 1994, chapter 16.4).



(a) (Two-sided) Bootstrap Hypothesis Testing based on empirical distribution from bootstrap samples.



(b) (Two-sided) Standard Hypothesis test with distributional assumption.

Figure 5: Bootstrap vs Distributional Assumption Hypothesis Testing.

Practically, we calculate estimates of the  $p$  values (where for  $p < \alpha$  we reject  $H_0$ ) as

$$\begin{aligned}\hat{p}_{\text{right}}^* &= \frac{1}{B} \sum_{b=1}^B \mathbb{1}_{[t^{(b)} \geq t_{\text{obs}}]} \\ \hat{p}_{\text{left}}^* &= \frac{1}{B} \sum_{b=1}^B \mathbb{1}_{[t^{(b)} \leq t_{\text{obs}}]} \\ \hat{p}_{\text{two-sided}}^* &= \frac{1}{B} \left[ \min \left\{ \sum_{b=1}^B \mathbb{1}_{[t^{(b)} \geq t_{\text{obs}}]}, \sum_{b=1}^B \mathbb{1}_{[t^{(b)} \leq t_{\text{obs}}]} \right\} \right]\end{aligned}\tag{8}$$

#### 1.4.1 Bootstrap Hypothesis Test on the Difference of Means I: Theory

Consider we have two samples  $X^{(1)} = (x_1^{(1)}, \dots, x_n^{(1)})$  and  $X^{(2)} = (x_1^{(2)}, \dots, x_n^{(2)})$ . And we want to test the null-hypothesis

$$H_0 : \mu^{(2)} - \mu^{(1)} \leq c \quad \text{vs} \quad H_1 : \mu^{(2)} - \mu^{(1)} > c\tag{9}$$

where  $c \in \mathbb{R}$ . Assume the true means are  $\mu^{(1)}$  and  $\mu^{(2)}$ . We can generate a bootstrap distribution under  $H_0$  by sampling from

$$\begin{aligned}z_i^{(1)} &= x_i^{(1)} - \bar{X}_{\text{sample}}^{(1)} - \frac{c}{2} \quad \text{and} \quad z_i^{(2)} = x_i^{(2)} - \bar{X}_{\text{sample}}^{(2)} + \frac{c}{2} \\ Z^{(1)} &= (z_1^{(1)}, \dots, z_n^{(1)}) \quad \text{and} \quad Z^{(2)} = (z_1^{(2)}, \dots, z_n^{(2)})\end{aligned}\tag{10}$$

and using the test statistic

$$\begin{aligned}t_{\text{obs}} &= \frac{\bar{X}_{\text{sample}}^{(2)} - \bar{X}_{\text{sample}}^{(1)} - c}{S_p \sqrt{\frac{1}{n^{(1)}} + \frac{1}{n^{(2)}}}} \\ S_p &= \sqrt{\frac{(n^{(1)} - 1) (S^{(1)})^2 + (n^{(2)} - 1) (S^{(2)})^2}{n^{(1)} + n^{(2)} - 2}}\end{aligned}\tag{11}$$

with analogous calculations of  $t^{(b)}$  on the bootstrap samples (the basic reasoning behind the denominator is that the scale on which means can be compared is the standard deviation, without which a statement like the means are 1 apart is meaningless (very different statement for standard deviations of the order of 10 or  $10^{-2}$ )).

### 1.4.2 Bootstrap Hypothesis Test on the Difference of Means I: Example Application

To test the described method, we use samples from two normal distributions where we specify the means and standard deviations. Based on our results, the test seems appropriate. Two examples are shown in figure 6. Note that as in hypothesis testing we are very cautious and go the devils-advocate route, we only reject  $H_0$  if we are very sure so in the test with large standard deviation (figure 6b) we wrongly accept  $H_0$ , which is still a reasonable choice given a standard deviation of 2 in the distributions and only a difference in means of 1.

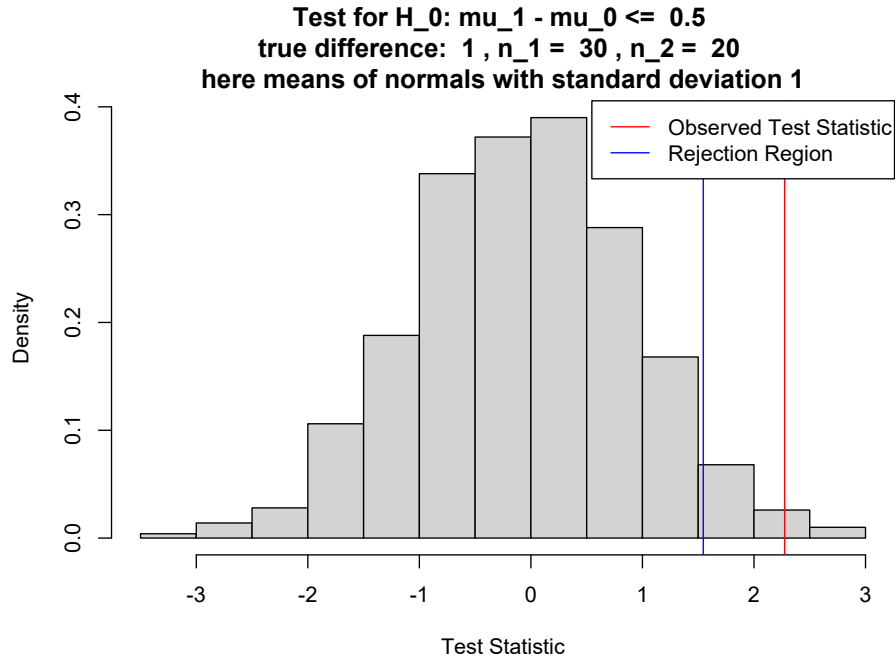
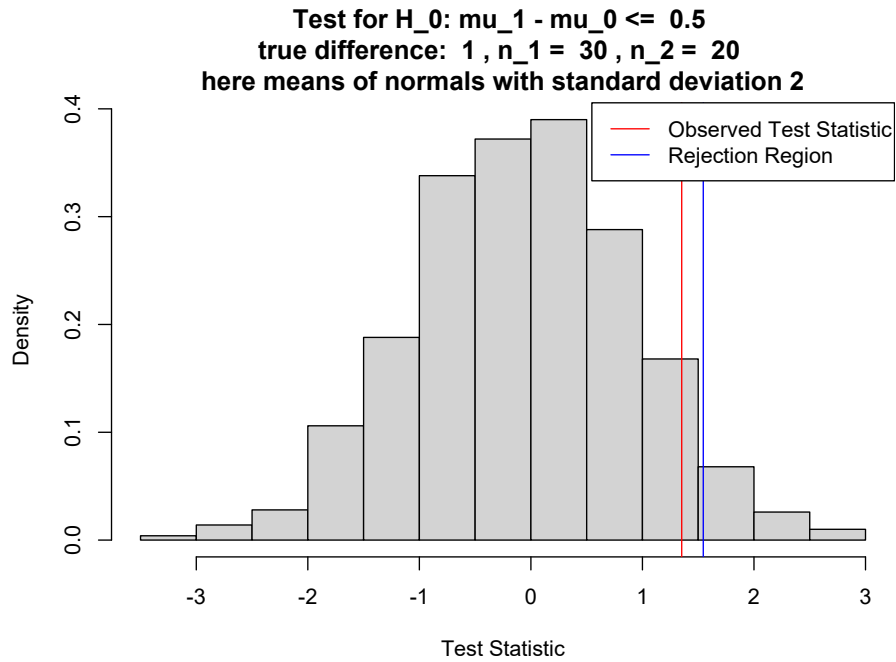
(a) Correctly Rejecting  $H_0$ .(b) Wrongly Accepting  $H_0$ .

Figure 6: Bootstrap Hypothesis Test on the Difference of Means

### 1.4.3 Bootstrap Hypothesis Test on the Difference of Means III: Task 1c

We choose a  $\alpha = 0.05$ -significance level and test  $H_0 : \mu^{(2)} - \mu^{(1)} \leq 0.25$  vs  $H_1 : \mu^{(2)} - \mu^{(1)} > 0.25$ .



Let us apply the previously introduced and tested methods to the samples from the exercise. Based on the samples we get a  $p$ -value of  $\hat{p}_{\text{right}}^* = 0.254$ , so  $H_0$  has to be accepted. As previously discussed, it would be unreasonable to reject  $H_0$  as of the large stray in the distributions compared to the difference in means (see figure 3 again). The test distribution is shown in figure 7 and the code to obtain the results in code-snippet 4.

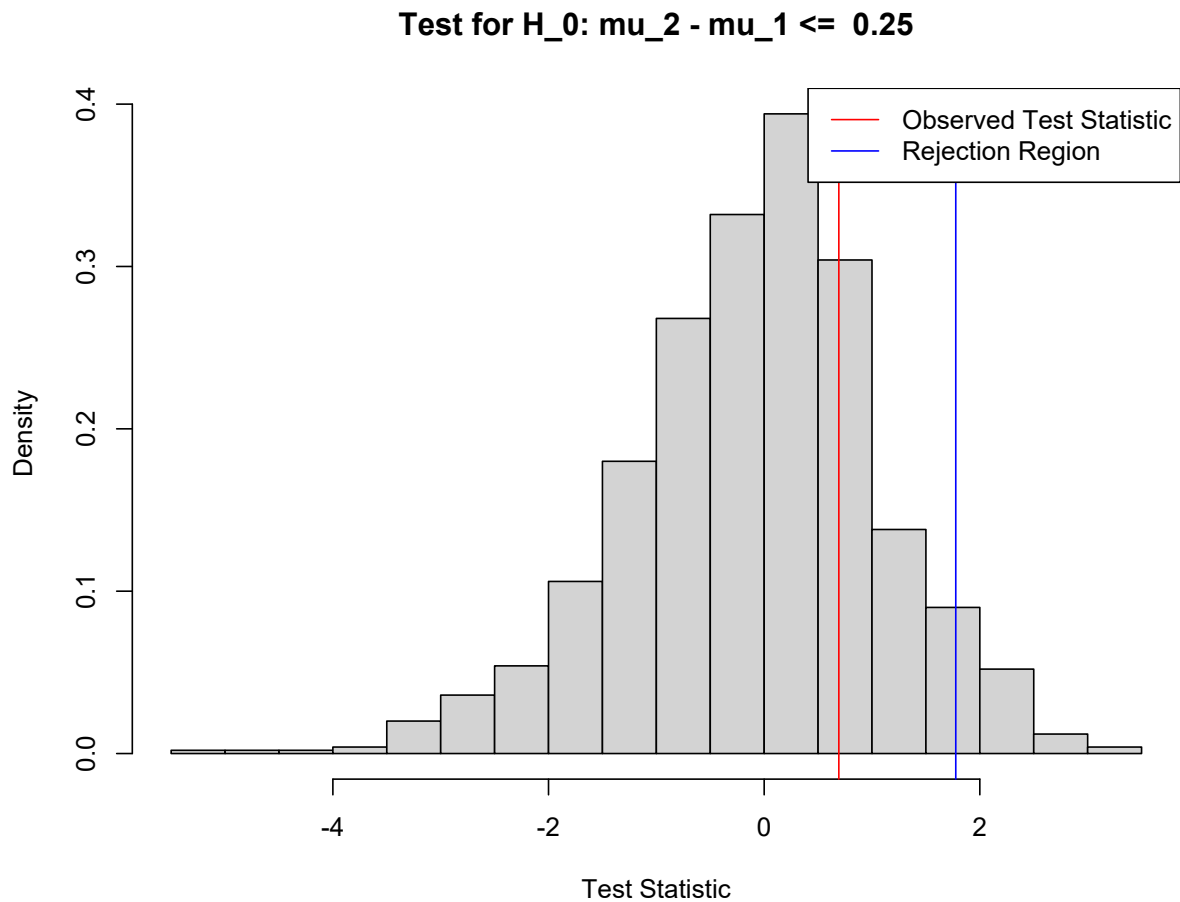


Figure 7: Bootstrap distribution of the test statistic  $t^{(b)}$  for the samples  $X^{(1)}$  and  $X^{(2)}$  for  $B = 1000$ .

```

1 bootstrap_mean_difference_greater_than_test <- function(sample1,
  ↪ sample2, diff, B = 1000) {
2   # sample1: vector of observations from sample 1
3   # sample2: vector of observations from sample 2
4   # diff: H_0: mu_1 - mu_2 <= diff
5   # B: number of bootstrap samples to use
6   # returns: evals of statistic on bootstrap samples, p-value for test
7
8   n1 <- length(sample1)
9   n2 <- length(sample2)
10
11   # test statistic used
12   t <- function(a, b) {
13     S <- sqrt(((n1 - 1) * var(a) + (n2 - 1) * var(b)) / (n1 + n2 -
  ↪ 2))
14     Tstat <- (mean(a) - mean(b) - diff) / (S * sqrt(1/n1 + 1/n2))
15     return(Tstat)
16   }
17
18   # observed test statistic
19   t_obs <- t(sample1, sample2)
20
21   # adapted samples for generating a bootstrap distribution under H_0
22   # note that we test for the extreme, so mu_1 - mu_2 = diff
23   # and then make a right-sided test
24
25   s1adapt <- sample1 - mean(sample1) + diff / 2
26   s2adapt <- sample2 - mean(sample2) - diff / 2
27
28   # generate bootstrap samples
29   tbs <- replicate(B, t(sample(s1adapt, n1, replace = TRUE),
  ↪ sample(s2adapt, n2, replace = TRUE)))
30
31   # compute p-value
32   pval <- mean(tbs >= t_obs)
33
34   return(list("t_obs" = t_obs, "tbs" = tbs, "pval" = pval))
35 }
36
37 # apply to the given samples, switch the order of the arguments, as
38 # I programmed the function before noticing, mu_1 - mu_2 is not the
39 # test the teacher probably had in mind
40 res <- bootstrap_mean_difference_greater_than_test(sample2, sample1,
  ↪ 0.25, 1000)

```

Code-Snippet 4: Bootstrap hypothesis test for  $H_0 : \mu^{(2)} - \mu^{(1)} \leq 0.25$  vs  $H_1 : \mu^{(2)} - \mu^{(1)} > 0.25$ .

## 1.5 Bootstrapping on Linear Regression Coefficients (task 2)

In this task, we regress on the waiting time between eruptions of a geyser given the length of the subsequent eruption.

### 1.5.1 Linear Fit on the Geysir Data (task 2a)

The code for the linear fit is given in code-snippet 5 and the fit is shown in figure 8. We get  $\beta_0 = 99 \pm 2$  and  $\beta_1 = -7.8 \pm 0.6$  (rounding the errors up and the coefficients to the significant digits). Note that based on the plot we can see that using a linear model only makes limited sense.

```
1      # Load the Data
2      library(MASS)
3      data(geyser)
4      # Fit the Model
5      model <- lm(waiting ~ duration, data = geyser)
6      # Summarize the Model
7      print(summary(model))
8      # Plot the Model
9      plot(geyser$duration, geyser$waiting, xlab = "Duration", ylab =
10     ↪ "Waiting", main = "Waiting vs. Duration Linear Regression")
11     abline(model, col = "red")
```

Code-Snippet 5: Linear fit on the geysir data.

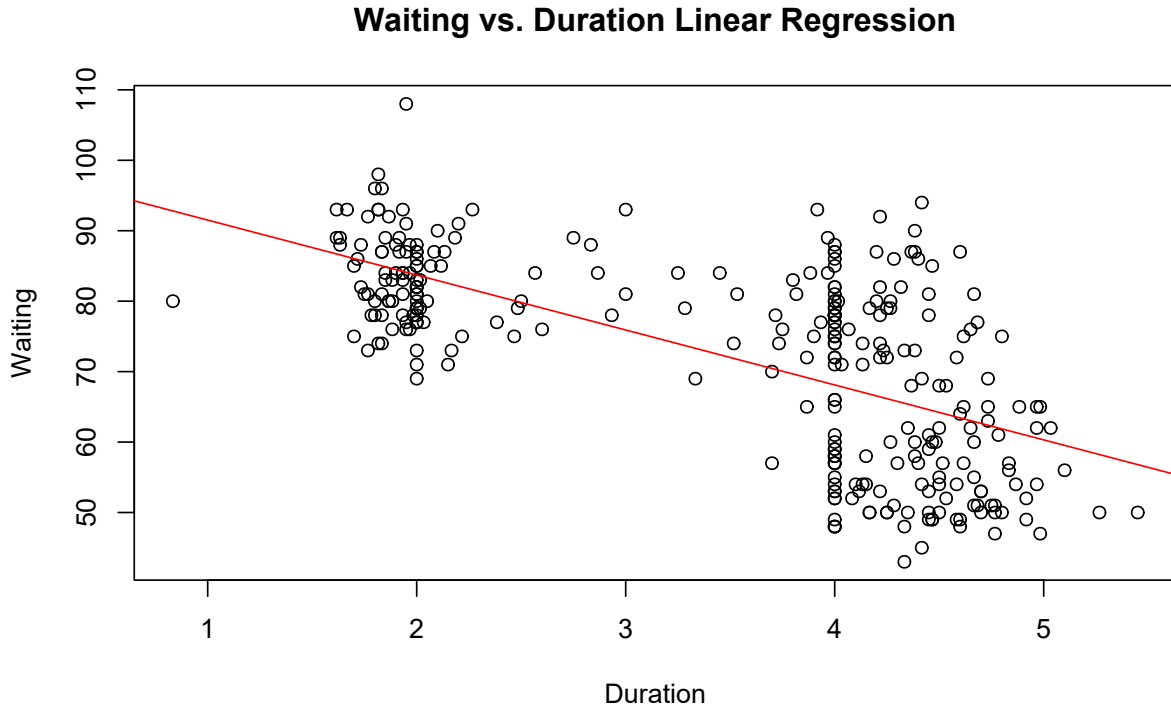


Figure 8: Linear fit on the geysir data.

## 1.6 Bootstrap estimates by resampling cases (task 2b)

By resampling data points and fitting multiple models, we can bootstrap-estimate the standard error of the coefficients. We just employ the previously introduced bootstrap methods, caveats mentioned there still apply.

The code for the bootstrap estimates is given in code-snippet 6. We get  $\hat{\beta}_0^* \approx 99.3$ ,  $\text{SE}(\hat{\beta}_0) \approx 1.4$ ,  $\hat{\beta}_1^* \approx -7.8$  and  $\text{SE}(\hat{\beta}_1) \approx 0.5$ .

```

1 bootstrap_lin_regr <- function(x, y, B = 1000) {
2   n <- length(x)
3   # point-estimate
4   mod <- lm(y ~ x)
5   b0_hat <- mod$coefficients[1]
6   b1_hat <- mod$coefficients[2]
7   # bootstrap
8   b0_bs <- rep(NA, B)
9   b1_bs <- rep(NA, B)
10  # bootstrap samples
11  for (i in 1:B) {
12    idx <- sample(1:n, n, replace = TRUE)
13    mod_b <- lm(y[idx] ~ x[idx])
14    b0_bs[i] <- mod_b$coefficients[1]
15    b1_bs[i] <- mod_b$coefficients[2]
16  }
17  # bootstrap estimates
18  b0_hat_star <- mean(b0_bs)
19  b1_hat_star <- mean(b1_bs)
20  # standard errors
21  se_b0 <- sd(b0_bs)
22  se_b1 <- sd(b1_bs)
23  # nicely format the output
24  return(list("b0_hat" = b0_hat, "b1_hat" = b1_hat, "b0_bs" =
    ↪ b0_bs, "b1_bs" = b1_bs, "b0_hat_star" = b0_hat_star,
    ↪ "b1_hat_star" = b1_hat_star, "se_b0" = se_b0, "se_b1" =
    ↪ se_b1))
25 }
26
27 est <- bootstrap_lin_regr(geyser$duration, geyser$waiting, B = 1000)
28 print(est)
29 # calculate the confidence intervals
30 # function to do so previously defined
31 alpha <- 0.04
32 b0_interval <- standard_normal_interval(est$se_b0, est$b0_hat, alpha
    ↪ = alpha)
33 b1_interval <- standard_normal_interval(est$se_b1, est$b1_hat, alpha
    ↪ = alpha)

```

Code-Snippet 6: Bootstrap estimates on the geysir data.

## 1.7 Standard normal bootstrap confidence intervals (task 2c)

Based on the code from code-snippet 2 we get the normal bootstrap confidence interval for  $\alpha = 0.04$  as  $[96.45, 102.16]$  for  $\beta_0$  and  $[-8.73, -6.87]$  for  $\beta_1$ .

## 1.8 Testing on $H_0 : \beta_0 = 0$ or rather $H_0 : \beta_0 = c$ (task 2d)

We are supposed to test  $H_0 : \beta_0 = 0$  vs  $H_1 : \beta_0 \neq 0$ . However, given the small standard error and the large value of  $\beta_0$ , p-values will be pretty much zero, so we will always reject  $H_0$ .

Let us therefore more generally test  $H_0 : \beta_0 = c$  vs  $H_1 : \beta_0 \neq c$  for  $c \in \mathbb{R}$  using the test statistic

$$T = \sqrt{\frac{nS_{xx}}{\sum_{i=1}^N}} \frac{\hat{\beta}_0 - c}{\hat{\sigma}} \sim t_{n-2} \quad (12)$$

$$S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2, \quad \hat{\sigma}^2 = \frac{\sum_{i=1}^N (Y_i - \bar{Y})^2 - \hat{\beta}_1^2 S_{xx}}{n-2}$$

### Test using the given test statistic (task 2d.i)

The code for doing the two-sided test is given in code-snippet 7. We choose an  $\alpha = 0.05$ -significance level and get report p-values for different values of  $c$  in table 3.

$c$	$p$ -value	reject $H_0 : \beta_0 = c$ at $\alpha = 0.05$
0	0	yes
96	9.1%	no
98	50.4%	no
100	72.5%	no
102	17.0%	no
104	1.7%	yes

Table 3: Hypothesis test for  $H_0 : \beta_0 = c$  vs  $H_1 : \beta_0 \neq c$  for  $c \in \mathbb{R}$ .

```

1  beta_0_statistic <- function(x, y, c, beta1, beta0) {
2    n <- length(x)
3    Sxx <- sum((x - mean(x))^2)
4    sigma <- sqrt((sum((y - mean(y))^2) - beta1^2 * Sxx) / (n - 2))
5    T <- sqrt((n * Sxx) / (sum(x^2))) * (beta0 - c) / sigma
6    return(T)
7  }
8
9  beta_0_t_test <- function(x, y, c, B = 1000) {
10    # calculate point estimates
11    mod <- lm(y ~ x)
12    beta0 <- mod$coefficients[1]
13    beta1 <- mod$coefficients[2]
14    # calculate the test statistic
15    t_obs <- beta_0_statistic(x, y, c, beta1, beta0)
16    n <- length(x)
17    # get p-value from t-distribution with n-2 degrees of freedom
18    p_value <- 2 * pt(-abs(t_obs), df = n - 2)
19    # nicely format the output
20    return(list("t_obs" = t_obs, "p_value" = p_value))
21  }
22
23  # do the t-test
24  res <- beta_0_t_test(geyser$duration, geyser$waiting, c = 104, B =
25    ↪ 1000)
26  print(res)

```

Code-Snippet 7: Hypothesis test for  $H_0 : \beta_0 = c$  vs  $H_1 : \beta_0 \neq c$  for  $c \in \mathbb{R}$ .

### Test using the bootstrap distribution (task 2d.ii)

Note that the estimate of  $\hat{\beta}_1$  does not change, when we mean-shift  $Y$ . We can therefore enforce  $H_0$  on the data, by

$$c = \bar{Y}^* - \hat{\beta}_1 \bar{x} \rightarrow \bar{Y}^* = \hat{\beta}_1 \bar{x} + c \rightarrow y_i^* = y_i - \bar{Y} + \bar{Y}^* = y_i - \bar{Y} + \hat{\beta}_1 \bar{x} + c \quad (13)$$

The code for doing the two-sided test is given in code-snippet 8. We choose an  $\alpha = 0.05$ -significance level and get report p-values for different values of  $c$  in table 4.

As expected the bootstrap distribution is more narrow than the  $t$ -distribution, so the bootstrap test rejects  $H_0$  more often than the  $t$ -test.

```

1  beta_0_statistic <- function(x, y, c, beta1, beta0) {
2    n <- length(x)
3    Sxx <- sum((x - mean(x))^2)
4    sigma <- sqrt((sum((y - mean(y))^2) - beta1^2 * Sxx) / (n - 2))
5    T <- sqrt((n * Sxx) / (sum(x^2))) * (beta0 - c) / sigma
6    return(T)
7  }
8
9  beta_0_bootstrap_test <- function(x, y, c, B = 1000) {
10     # calculate point estimates
11     mod <- lm(y ~ x)
12     beta0 <- mod$coefficients[1]
13     beta1 <- mod$coefficients[2]
14     # calculate the test statistic
15     t_obs <- beta_0_statistic(x, y, c, beta1, beta0)
16     # impose the null hypothesis on the data
17     y <- y - mean(y) + beta1 * mean(x) + c
18     # check that we correctly imposed the null hypothesis
19     beta0new <- lm(y ~ x)$coefficients[1]
20     if (abs(beta0new - c) > 1e-10) {
21       stop("Error: Null hypothesis not imposed correctly")
22     }
23     # bootstrap
24     t_bs <- rep(0, B)
25     for (i in 1:B) {
26       idx <- sample(1:length(x), length(x), replace = TRUE)
27       # for H0 we assume that beta1 is known
28       # if we would estimate beta1 from the data
29       # the null hypothesis would not be strictly imposed
30       t_bs[i] <- beta_0_statistic(x[idx], y[idx], c, beta1,
31         ↪ lm(y[idx] ~ x[idx])$coefficients[1])
32     }
33     # calculate the p-value (two-sided)
34     p_left <- mean(t_bs <= t_obs)
35     p_right <- mean(t_bs >= t_obs)
36     p_value <- min(p_left, p_right) * 2
37     # nicely format the output
38     return(list("t_obs" = t_obs, "t_bs" = t_bs, "p_value" =
39       ↪ p_value))
40   }
41
42   # do the bootstrap test
43   res <- beta_0_bootstrap_test(geyser$duration, geyser$waiting, c =
44     ↪ 99, B = 1000)

```

Code-Snippet 8: Bootstrap hypothesis test for  $H_0 : \beta_0 = c$  vs  $H_1 : \beta_0 \neq c$  for  $c \in \mathbb{R}$ .



$c$	$p$ -value	reject $H_0 : \beta_0 = c$ at $\alpha = 0.05$
0	0	yes
96	1.4%	yes
98	35%	no
100	65.8%	no
102	8.4%	no
104	0.4%	yes

Table 4: Bootstrap hypothesis test for  $H_0 : \beta_0 = c$  vs  $H_1 : \beta_0 \neq c$  for  $c \in \mathbb{R}$ .

## 2 Jackknife and Further Bootstrap Applications

Consider the following sample from  $X \sim \mathcal{N}(\mu, \sigma^2)$ :

$$X_{\text{sample}} = \{2.8, 3, 2.3, 1.3, 1.4, 3.2, 2.9, 1.1, 2.2, 3.5, 1.8, 3.5, 2.5, 2, 2.8\} \quad (14)$$

### 2.1 Bias correction and standard error estimation using Bootstrap and Jackknife (task 3a, 3b)

Consider the estimator of the mean

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i \quad (15)$$

Using the Bootstrap method (as introduced in section 1), we get the results shown in table 5 using the code in code-snippet 9 (task b). As the calculated bias is smaller than the standard error, it makes no sense to correct by it. This is also the right conclusion, as the mean is an unbiased estimator.

point estimate $\hat{\mu}$	$\hat{\text{SE}}(\hat{\mu})$	bootstrap estimate $\hat{\mu}^*$	$\hat{\text{bias}}(\hat{\mu})$	bias-corrected estimate $\bar{\mu}$
$\approx 2.420$	$\approx 2.415$	$\approx 0.19$	$\approx -0.0053$	$\approx 2.4253$

Table 5: Bootstrap estimates regarding  $\mu$  for  $X_{\text{sample}}$  for  $B = 1000$ .

Different from bootstrapping, in the Jackknife method, we generate  $n$  subsamples by removing one observation at a time,

$$X_{\text{sample}(-i)} = \{x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n\} \quad (16)$$

which is computationally less expensive than bootstrapping.

**Jackknife-Intuition:** Note, that the idea of the Jackknife is a bit different. In Jackknife,

```

1 bootstrap <- function(X, T, B = 1000) {
2   # Sample estimate
3   sample_estimate <- T(X)
4   n <- length(X)
5
6   # Generate bootstrap samples and apply function T
7   bootstrap_samples <- replicate(B, T(sample(X, size = n, replace
8     ↪ = TRUE)))
9
10  # Sort the bootstrap samples
11  stat <- sort(bootstrap_samples)
12
13  # The estimator is the mean
14  bootstrap_estimate <- mean(stat)
15
16  # The standard error is the standard deviation of the statistic
17  sd <- sd(stat)
18
19  # The estimated bias is the difference between the sample
20  ↪ estimate
21  # and the bootstrap estimate
22  bias <- bootstrap_estimate - sample_estimate
23
24  # Bias-corrected estimate
25  corrected_estimate <- sample_estimate - bias
26
27  # Return the estimate, standard error, and bias
28  return(list(sample_estimate = sample_estimate,
29    bootstrap_estimate = bootstrap_estimate,
30    sd = sd,
31    bias = bias,
32    corrected_estimate = corrected_estimate))
33 }
34 X <- c(2.8, 3, 2.3, 1.3, 1.4, 3.2, 2.9, 1.1, 2.2, 3.5, 1.8, 3.5,
35   ↪ 2.5, 2, 2.8)
36 # Bootstrap estimate of the mean
37 res <- bootstrap(X, mean)
38 print(res)

```

Code-Snippet 9: Code for Bootstrap estimates regarding  $X_{\text{sample}}$ .

we ask ourselves, what happens, when one subsample to the other differs by having one observation each the other does not have while sharing the rest (so similar subsamples). For the estimate of the bias and variance of the estimator we then have to scale this up to "what if we had totally different samples". This leads to the Jackknife estimates

$$\hat{t}_{\text{jack}} = \frac{1}{n} \sum_{i=1}^n \hat{t}_{(-i)}, \quad \hat{\text{bias}}_{\text{jack}} = (n-1) (\hat{t}_{\text{jack}} - \hat{t}), \quad \hat{\text{SE}}_{\text{jack}} = \sqrt{\frac{n-1}{n} \sum_{i=1}^n (\hat{t}_{(-i)} - \hat{t}_{\text{jack}})^2} \quad (17)$$

Note that for the case of the mean

$$\begin{aligned} \hat{\mu}_{\text{jack}} &= \frac{1}{n} \sum_{i=1}^n \hat{\mu}_{(-i)} \\ &= \frac{1}{n} \sum_{i=1}^n \frac{1}{n-1} \sum_{j=1, j \neq i}^n x_j \\ &= \frac{1}{n} \sum_{i=1}^n \left( \frac{1}{n-1} \left( \sum_{j=1}^n x_j - x_i \right) \right) \\ &= \frac{1}{n} \sum_{i=1}^n \left( \frac{1}{n-1} (n\hat{\mu} - x_i) \right) \\ &= \frac{n}{n-1} \hat{\mu} - \frac{1}{(n-1)} \hat{\mu} \\ &= \hat{\mu} \end{aligned} \quad (18)$$

so the Jackknife estimate of the mean is actually just the sample mean itself (at least from my calculations), so we do not expect any bias so also no bias correction. The numerical results are shown in table 6 using the code in code-snippet 10 (task a). As expected, the Jackknife estimate of the mean is the sample mean itself, so we reflect, that there is no bias and nothing to correct.

point estimate $\hat{\mu}$	$\hat{\text{SE}}(\hat{\mu})$	jackknife estimate $\hat{\mu}_{\text{jack}}$	$\hat{\text{bias}}(\hat{\mu})$
$\approx 2.42$	$\approx 0.20$	$\approx 2.42$	$= 0$

Table 6: Jackknife estimates regarding  $\mu$  for  $X_{\text{sample}}$ .

```

1  jackknife <- function(X, T) {
2      n <- length(X)
3
4      # Array to store jackknife samples
5      jackknife_samples <- numeric(n)
6
7      # Jackknife samples, leaving one out
8      for (i in 1:n) {
9          jackknife_samples[i] <- T(X[-i])
10     }
11
12     # The sample estimate
13     sample_estimate <- T(X)
14
15     # The jackknife estimate is the mean of the jackknife samples
16     jackknife_estimate <- mean(jackknife_samples)
17
18     # Standard error calculation
19     jackknife_se <- sqrt((n - 1) / n * sum((jackknife_samples -
20     ↪ jackknife_estimate)^2))
21
22     # Bias estimation
23     bias <- (n - 1) * (jackknife_estimate - sample_estimate)
24
25     # Bias-corrected estimate
26     corrected_estimate <- sample_estimate - bias
27
28     # Return the estimate, standard error, and bias
29     return(list(sample_estimate = sample_estimate,
30                jackknife_estimate = jackknife_estimate,
31                se = jackknife_se,
32                bias = bias,
33                corrected_estimate = corrected_estimate))
34 }
35
36 # sample
37 X <- c(2.8, 3, 2.3, 1.3, 1.4, 3.2, 2.9, 1.1, 2.2, 3.5, 1.8, 3.5,
38 ↪ 2.5, 2, 2.8)
39
40 # Jackknife estimate of the mean
41 res <- jackknife(X, mean)
42 print(res)

```

Code-Snippet 10: Code for Jackknife estimates regarding  $X_{\text{sample}}$ .

## 2.2 Non-parametric bootstrap method to test the hypothesis $H_0 : X \sim \mathcal{N}(2, 1)$ (task 3c)

Let us first get a feeling for how realistic  $H_0$  might be, based on the plot in figure 9. The centering seems to be somewhat off, but just from looking at the data, I could not straight reject  $H_0$ . The best we could do, I would say, is to give the likelihood, that this data comes from  $\mathcal{N}(2, 1)$ , but let us do a fancy test with Kolmogorov-Smirnov statistic instead<sup>3</sup>.

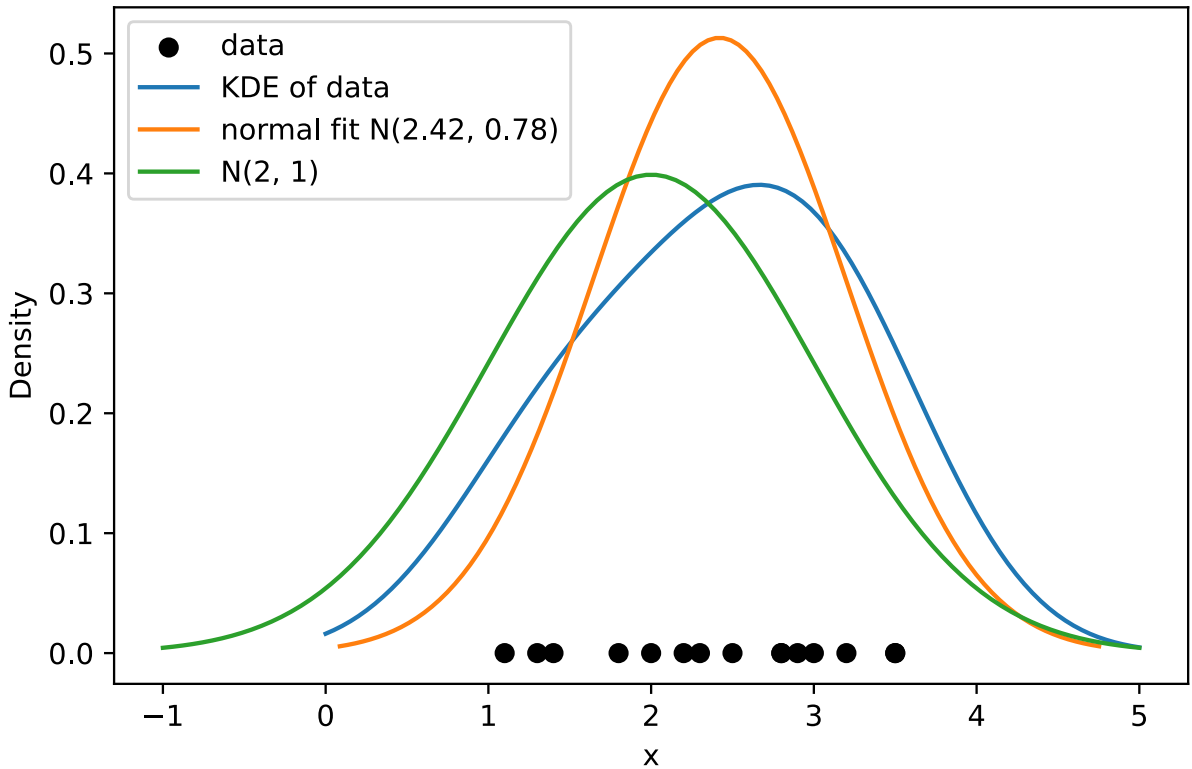


Figure 9: Kernel density estimate of  $X_{\text{sample}}$  with normal fit.

The Kolmogorov-Smirnov statistic is defined as

$$D = \max_{i=1, \dots, n} \left\{ \frac{i}{n} - F_0(x_{(i)}), F_0(x_{(i)}) - \frac{i-1}{n} \right\}, \quad F_0(x) = P_0(X \leq x) \text{ under } H_0 \quad (19)$$

As in the task it is given, that the data comes from some normal distribution, we can impose  $H_0$  onto the empirical distribution by sampling from

$$x_i^* = \frac{x_i - \bar{X}}{\sigma_X} \cdot \sigma_0 + \mu_0, \quad \sigma_0 = 1, \quad \mu_0 = 2 \quad (20)$$

<sup>3</sup>Testing has the very inherent problem that we try to crunch the information of data into a single yes/no decision with a p-value, which one who lacks integrity might p-hack on.

---

Using the code in code-snippet 11, we get a p-value of 74.2%, so we cannot reject  $H_0$  at a significance level of 5%. Just for comparison, the p-value for the test for  $H_0 : X \sim \mathcal{N}(1.5, 1)$  is 3.4%, so we can reject  $H_0$  at a significance level of 5%.

```

1  # to have pnorm available
2  library(stats)
3  # Kolmogorov-Smirnov statistic
4  ks_stat <- function(sample, F0) {
5      n <- length(sample)
6      sorted_sample <- sort(sample)
7      # initialize vector of KS values
8      ks_vals <- rep(NA, n)
9      for (i in 1:n) {
10         ks_vals[i] <- max(i/n - F0(sorted_sample[i]),
11             ↪ F0(sorted_sample[i]) - (i-1)/n)
12     }
13     # return the maximum KS value
14     return(max(ks_vals))
15 }
16 # Bootstrap test
17 bootstrap_ks_test <- function(X, F0, sd0, mu0, B = 1000) {
18     n <- length(X)
19     original_ks <- ks_stat(X, F0)
20     bootstrap_ks <- rep(NA, n)
21     # impose H_0 on the data, assuming
22     # X is normally distributed
23     X <- (X - mean(X)) / sd(X) * sd0 + mu0
24     for (b in 1:B) {
25         # Generate bootstrap sample and compute KS stat
26         sample <- sample(X, size = n, replace = TRUE)
27         bootstrap_ks[b] <- ks_stat(sample, F0)
28     }
29     # Compute p-value
30     p_left <- mean(bootstrap_ks <= original_ks)
31     p_right <- mean(bootstrap_ks >= original_ks)
32     p_value <- min(p_left, p_right) * 2
33     list(original_ks = original_ks, p_value = p_value)
34 }
35 # Sample data X
36 X <- c(2.8, 3, 2.3, 1.3, 1.4, 3.2, 2.9, 1.1, 2.2, 3.5, 1.8, 3.5,
37     ↪ 2.5, 2, 2.8)
38 # Null hypothesis parameters
39 mu0 <- 2
40 sd0 <- 1
41 # Null hypothesis:  $X \sim N(\mu_0, \sigma_0)$ 
42 F0 <- function(x) pnorm(x, mean = mu0, sd = sd0)
43 # Perform the bootstrap KS test
44 result <- bootstrap_ks_test(X, F0, sd0, mu0, B = 1000)
45 print(result)

```

Code-Snippet 11: Code for Kolmogorov-Smirnov test regarding  $X_{\text{sample}}$ .

# References

- Efron, Bradley and R.J. Tibshirani (May 1994). *An Introduction to the Bootstrap*. Chapman and Hall/CRC. DOI: 10.1201/9780429246593. URL: <http://dx.doi.org/10.1201/9780429246593>.
- Hesterberg, Tim C. (2015). »What Teachers Should Know About the Bootstrap: Resampling in the Undergraduate Statistics Curriculum«. In: *The American Statistician* 69.4. PMID: 27019512, pages 371–386. DOI: 10.1080/00031305.2015.1089789. URL: <https://doi.org/10.1080/00031305.2015.1089789>.
- Overduin, Stephen (2004). »USE OF THE LOGNORMAL DISTRIBUTION FOR SURVIVAL DATA: INFERENCE AND ROBUSTNESS«. Master's thesis. Simon Fraser University. URL: <https://core.ac.uk/download/pdf/56374045.pdf>.