

# Exercises Lecture 09

## Resampling 2: The Bootstrap

*SCR team*

*14 November, 2019*

### Exercises part 1

When the parametric bootstrap would not be an optimal procedure... but still could be used.

#### 1.1 Check the slides for the Parametric Bootstrap of Reaction Time.

We have created our own simulated reaction time data-set again, but now we use  $n = 1000$  reaction times:

```
set.seed(160945)
alpha <- 3; beta <- 1 # true values
n <- 1000
X <- rgamma(n, alpha, beta) # data
```

The Gamma distribution with parameters shape =  $\alpha$  and rate =  $\beta$ , and has density

$$f(x) = \frac{\beta^\alpha}{\Gamma(\alpha)} x^{(\alpha-1)} e^{-(\beta x)}$$

the mean and variance are  $E(X) = \alpha/\beta$  and  $Var(X) = \alpha/\beta^2$ .

**a**

Use the function `dgamma()` and code your own density function `my_dgamma()` that gives the same results as `dgamma()`. Note that very small differences between the functions (e.g. a difference of  $1e - 15$ ) are allowed.

**Answer**

```
my_dgamma <- function(x, alpha, beta) {
  out_dens <- beta^alpha * x^(alpha - 1) * exp(-x/beta) / gamma(alpha)
  return(out_dens)
}
all.equal(dgamma(X, alpha, beta), my_dgamma(X, alpha, beta))
```

```
## [1] TRUE
```

**b**

You may have learned in your Probability & Statistics course(s) that the maximum likelihood estimates of the parameters of the gamma distribution for our specific data would be

```
beta_hat <- mean(X) / var(X)
alpha_hat <- mean(X) * beta_hat
c(shape = alpha_hat, rate = beta_hat)
```

```
##      shape      rate
## 3.037211 1.057207
```

In the slides of the lecture we have seen a parametric bootstrapping procedure to obtain an estimate of the standard error of the sampling distribution of the median. We used the maximum likelihood estimates of the shape and rate parameters to be able to draw samples from the gamma distribution (with the function `rgamma()`). What is your estimate of the expected median? What is your estimate of the standard error of the median?

**Answer:**

```
B <- 1e3; t_pboots <- numeric(B)
for (b in 1:B) {
  X_b <- rgamma(n, alpha_hat, beta_hat)
  t_pboots[b] <- median(X_b)
}
c("E[median]" = mean(t_pboots), "S_median." = sd(t_pboots))
```

```
## E[median] S_median.
## 2.56645941 0.06082013
```

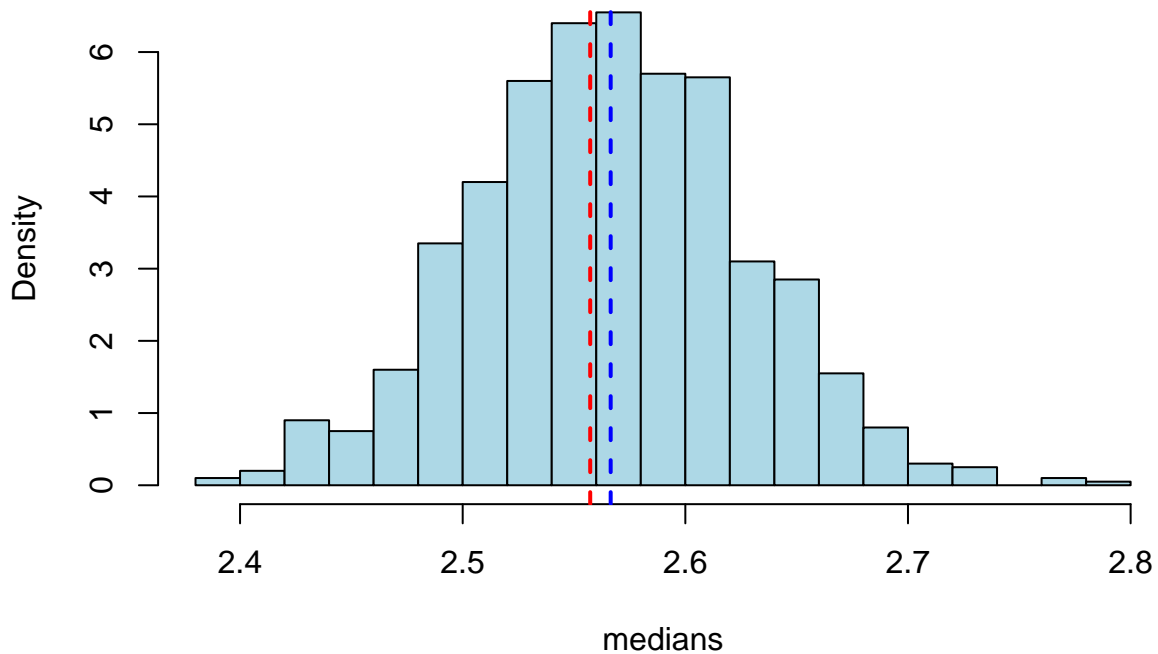
**c.**

Create a histogram of the parametric bootstrapped replicates of the median to visualize the estimate of the sampling distribution of the median. Add a vertical line of your observed median (in red), as well as the expected median (in blue) in your estimate of the sampling distribution of the median.

**Answer:**

```
hist(t_pboots, xlab = "medians",
     breaks = "FD", col = "lightblue",
     main = "Parametric Boot Samp. Distr.",
     freq = FALSE
)
abline(v = median(X), lwd = 2, lty = 2, col = "red")
abline(v = mean(t_pboots), lwd = 2, lty = 2, col = "blue")
```

## Parametric Boot Samp. Distr.



d.

Perform a Monte-Carlo study, by sampling  $B = 1000$  estimates of the median, but make sure to use the true shape ( $\alpha = 3$ ) and rate parameter ( $\beta = 1$ ).

d.i

Is the expectation of the Monte-Carlo medians the same as the expectation of the parametric bootstrapped medians? What is the difference (= estimate of the bias)?

d.ii

How about the standard error of the Monte-Carlo median vs. the parametric bootstrapped median? What would be your estimate of the bias?

d.iii

Visualize your Monte-Carlo replicates of the median in a histogram and add your observed median as a red vertical line and the expected median as a blue vertical line.

### Answer

The MC-study:

```
B <- 1e3; t_mcs <- numeric(B)
for (b in 1:B) {
  X_b <- rgamma(n, alpha, beta)
  t_mcs[b] <- median(X_b)
}
```

The expected median from the Monte Carlo experiment (`Emedian_mc`) is different from the expected median of the parametric bootstrapped median (`Emedian_pb`).

```
c(Emedian_mc = mean(t_mcs), Emedian_pb = mean(t_pboots), Bias_hat = mean(t_pboots) - mean(t_mcs))
```

```
## Emedian_mc Emedian_pb   Bias_hat
##  2.6722817  2.5664594 -0.1058223
```

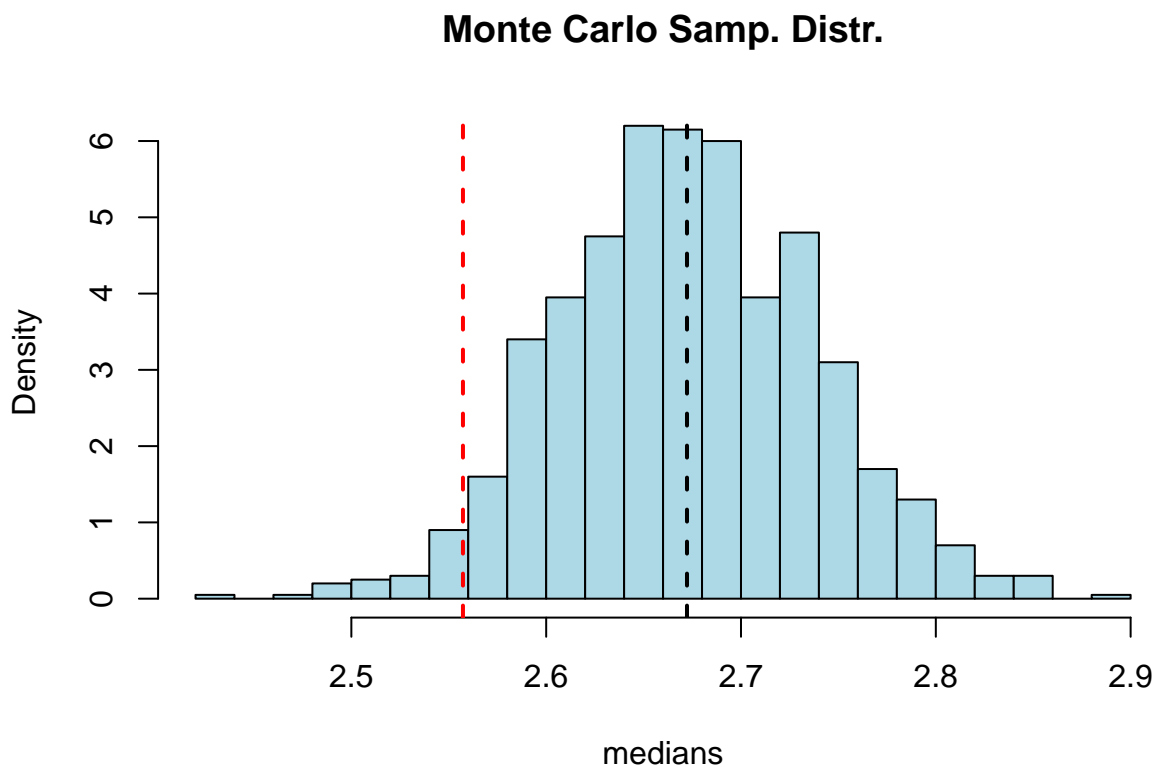
However, compared to the slides, the two standard error estimates of the median seem to be more similar to each other:

```
c(s_median_mc = sd(t_mcs), s_median_pb = sd(t_pboots), bias_se = abs(sd(t_mcs) - sd(t_pboots)))
```

```
## s_median_mc s_median_pb   bias_se
## 0.065367340 0.060820131 0.004547209
```

The histogram:

```
hist(t_mcs, xlab = "medians",
      breaks = "FD", col = "lightblue",
      main = "Monte Carlo Samp. Distr.",
      freq = FALSE
)
abline(v = median(X), lwd = 2, lty = 2, col = 'red')
# abline(v = mean(t_pboots), lwd = 2, lty = 2, col = "blue")
abline(v = mean(t_mcs), lwd = 2, lty = 2, col = "black")
```



e

Without conducting your experiment in R code, what would be a good way to get an idea of how ‘trustworthy’ all these estimates are? What would happen to your estimates for larger  $n$ , and what would happen for larger  $B$ ?

**Answer:**

Rerunning the whole experiment a number of times to get an idea of your simulation error.

For larger  $n$  the parametric bootstrap becomes more and more equal to the estimates of the Monte Carlo study (related to consistency).

For larger  $n$  the estimators would become more precise (efficient). A similar reasoning holds for larger  $B$ , since with larger  $B$  we would just reduce the simulation error.

## 1.2 A Parametric Bootstrap for Regression Analysis

Load the data set `Advertising.csv` either from

<http://www-bcf.usc.edu/~gareth/ISL/Advertising.csv>

or from this Rproject’s directory (path is `./data/Advertising.csv`) into *R* and explore the data set a bit.

This dataset (see James, Witten, Hastie and Tibshirani, 2017) contains measurements of sales (in thousands of units), and of TV, radio and newspaper budgets (in thousands of dollars), for 200 different markets.

Let  $y_i$  be the sales of a particular market  $i$ , let  $x_{i1}$  be the budget expenditure on TV advertizements in market  $i$ , let  $x_{i2}$  be the radio budget expenditure in market  $i$ , and let  $x_{i3}$  be the budget expenditure on newspapers in market  $i$ .

In this exercises we are interested in the linear model

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i3} + \epsilon_i,$$

where  $\beta_0$  is an intercept, and  $\beta_1, \beta_2, \beta_3$  are linear effects, and  $\epsilon_i$  is an “error” term for which we have the assumption

$$\epsilon_i \sim N(0, \sigma)$$

a)

Use `lm()` to perform a linear regression analysis of `sales` on `TV`, `radio` and `newspaper` according to the linear model described above. Store the results from your linear regression analysis into a variable (aka object), and take a look at the summary (`summary()`) of your results. Is there a significant contribution of “advertisement in newspapers” on sales in this data set (controlled for TV and radio advertisements)?

**Answer:**

```
adv_dat <- read.csv("./0_data/Advertising.csv")
adv_lm <- lm(sales ~ TV + radio + newspaper, data = adv_dat)
summary(adv_lm)

##
## Call:
## lm(formula = sales ~ TV + radio + newspaper, data = adv_dat)
```

```
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -8.8277 -0.8908  0.2418  1.1893  2.8292
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  2.938889   0.311908   9.422  <2e-16 ***
## TV           0.045765   0.001395  32.809  <2e-16 ***
## radio        0.188530   0.008611  21.893  <2e-16 ***
## newspaper   -0.001037   0.005871  -0.177    0.86
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.686 on 196 degrees of freedom
## Multiple R-squared:  0.8972, Adjusted R-squared:  0.8956
## F-statistic: 570.3 on 3 and 196 DF,  p-value: < 2.2e-16
```

Nope.

b)

The sales that only contains the intercept and the linear effect of newspaper advertizements can be denoted (and defined) as follows:

$$\hat{z}_i := \hat{\beta}_0 + \hat{\beta}_3 x_{i3} = y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \hat{\beta}_2 x_{i2} + \hat{\beta}_3 x_{i3} + \hat{\epsilon}_i),$$

Using the results that you have obtained in **a** create your own vector  $\hat{\mathbf{z}}$  that contains each  $\hat{z}_i$  for  $i \in 1 \dots 200$ .

**Answer:**

```
hat_betas <- coef(adv_lm)
hatsales_remain <- hat_betas[1] + hat_betas[4] * adv_dat$newspaper
```

or a bit more tedious;

```
hat_residuals <- residuals(adv_lm)
hatsales_remain2 <- {
  adv_dat$sales -
  hat_betas[2] * adv_dat$TV -
  hat_betas[3] * adv_dat$radio -
  hat_residuals
}
```

Just checking:

```
names(hatsales_remain2) <- names(hatsales_remain)
all.equal(hatsales_remain, hatsales_remain2)
```

```
## [1] TRUE
```

c)

Use `all.equal()` to verify that the estimated standard deviation of the errors (misnamed as “residual standard error” in the output of `summary.lm()`) can be computed as

$$\hat{\sigma} = \sqrt{\frac{\sum \hat{\epsilon}_i^2}{n - p}}$$

where  $n = 200$  markets, and  $p = 4$  (the number of parameters: intercept + linear effects).

**Answer**

```
hat_sigma <- sigma(adv_lm)
n <- 200; p <- 4
myhat_sigma <- sqrt(sum(residuals(adv_lm)^2) / (n - p))
all.equal(hat_sigma, myhat_sigma)
```

```
## [1] TRUE
```

c)

Suppose you are a trainee at an advertisement company, and your boss wants to stop the advertising in the newspaper. To be sure, he would like you to validate the estimated standard error of the linear effect  $\hat{\beta}_3$  (= the contribution in sales the newspaper).

Do this as follows: Create  $B = 1000$  parametric bootstrap samples of the  $\mathbf{z}$ , defined as

$$z_i^b = \hat{z}_i + \epsilon_i^b$$

where  $\epsilon_i^b$  is your own sampled residual from  $N(0, \hat{\sigma})$ , the normal distribution with mean zero and a standard deviation equal to your estimate of the variance of the errors ( $\hat{\sigma}^2$ ).

Then regress each  $\mathbf{z}^b$  (for  $b = 1 \dots B$ ) on the original `newspaper` variable ( $\mathbf{x}_3$ ), and save the estimate of the linear effect (the coefficient) of advertisement expenditure for the newspaper.

What is the mean of these bootstrapped main effects, and what is the standard deviation? Are these results similar to those obtained in a)?.

**Answer:**

```
B <- 1e3; beta_pboots <- rep(NA, B)
n <- nrow(adv_dat)
for (b in 1:B) {
  y <- hatsales_remain + rnorm(n, sd = hat_sigma)
  beta_pboots[b] <- coefficients(lm(y~adv_dat$newspaper))[2]
}
c(Estimate = mean(beta_pboots), "Std. Error" = sd(beta_pboots))
```

```
##      Estimate      Std. Error
## -0.0008270786  0.0055903864
```

```
summary(adv_lm)$coefficients[4,1:2]
```

```
##      Estimate   Std. Error
## -0.001037493  0.005871010
```

In general, for this specific example, the standard error as well as the mean estimate for the linear effect of `newspaper` from the parametric bootstrap procedure are very similar to the results obtained in the linear regression analysis.

## Exercises part 2

### 2.1 Combining the Empirical and Parametric Bootstrap for Regression Analysis

Repeat exercise 1.2d, but instead of sampling the residuals from a normal distribution with variance equal to the observed estimated residual variance, apply the empirical bootstrap. Do the results remain similar?

**Answer:**

```
B <- 1e3; beta_eboots <- rep(NA, B)
n <- nrow(adv_dat)
for (b in 1:B) {
  y <- hatsales_remain + sample(hat_residuals, n, replace = TRUE)
  beta_eboots[b] <- coefficients(lm(y~adv_dat$newspaper))[2]
}
c(Estimate = mean(beta_eboots), "Std. Error" = sd(beta_eboots))
```

```
##      Estimate   Std. Error
## -0.0008286854  0.0054799015
```

```
summary(adv_lm)$coefficients
```

```
##              Estimate Std. Error   t value    Pr(>|t|)
## (Intercept)  2.938889369 0.311908236  9.4222884 1.267295e-17
## TV           0.045764645 0.001394897 32.8086244 1.509960e-81
## radio        0.188530017 0.008611234 21.8934961 1.505339e-54
## newspaper    -0.001037493 0.005871010 -0.1767146 8.599151e-01
```

Yes, the results remain similar.

### 2.2 Failing the Bootstrap: Regression Analysis with too small $n$

Suppose we have the regression model

```
fm1 <- lm(Employed ~ ., data = longley)
M1 <- model.matrix(fm1)
betas <- coefficients(fm1)
summary(fm1)
```



```
##
## Call:
## lm(formula = Employed ~ ., data = longley)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.41011 -0.15767 -0.02816  0.10155  0.45539
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -3.482e+03  8.904e+02  -3.911 0.003560 **
## GNP.deflator  1.506e-02  8.492e-02   0.177 0.863141
## GNP          -3.582e-02  3.349e-02  -1.070 0.312681
## Unemployed   -2.020e-02  4.884e-03  -4.136 0.002535 **
## Armed.Forces -1.033e-02  2.143e-03  -4.822 0.000944 ***
## Population   -5.110e-02  2.261e-01  -0.226 0.826212
## Year          1.829e+00  4.555e-01   4.016 0.003037 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.3049 on 9 degrees of freedom
## Multiple R-squared:  0.9955, Adjusted R-squared:  0.9925
## F-statistic: 330.3 on 6 and 9 DF,  p-value: 4.984e-10
```

Program your own empirical bootstrap procedure on the residuals to see whether you can validate the standard errors of all coefficients (except the intercept) for the results in full model `fm1`.

**Answer:**

```
resids <- longley$Employed - fm1$fitted # step 2
B <- 1e3 ; set.seed(160913+43)
beta_stars <- t(apply(1:B, function(b){
  y <- M1 %*% betas + sample(resids, length(resids), replace = TRUE)
  bstar <- coefficients(lm(y ~ M1 - 1, data = longley))
  return(t(bstar))
}))
se_boot <- apply(beta_stars, 2, sd) # Bootstrap se's of the coefficients
se_boot
```

```
## [1] 6.602795e+02 6.403123e-02 2.523158e-02 3.677250e-03 1.596366e-03
## [6] 1.683030e-01 3.374880e-01
```

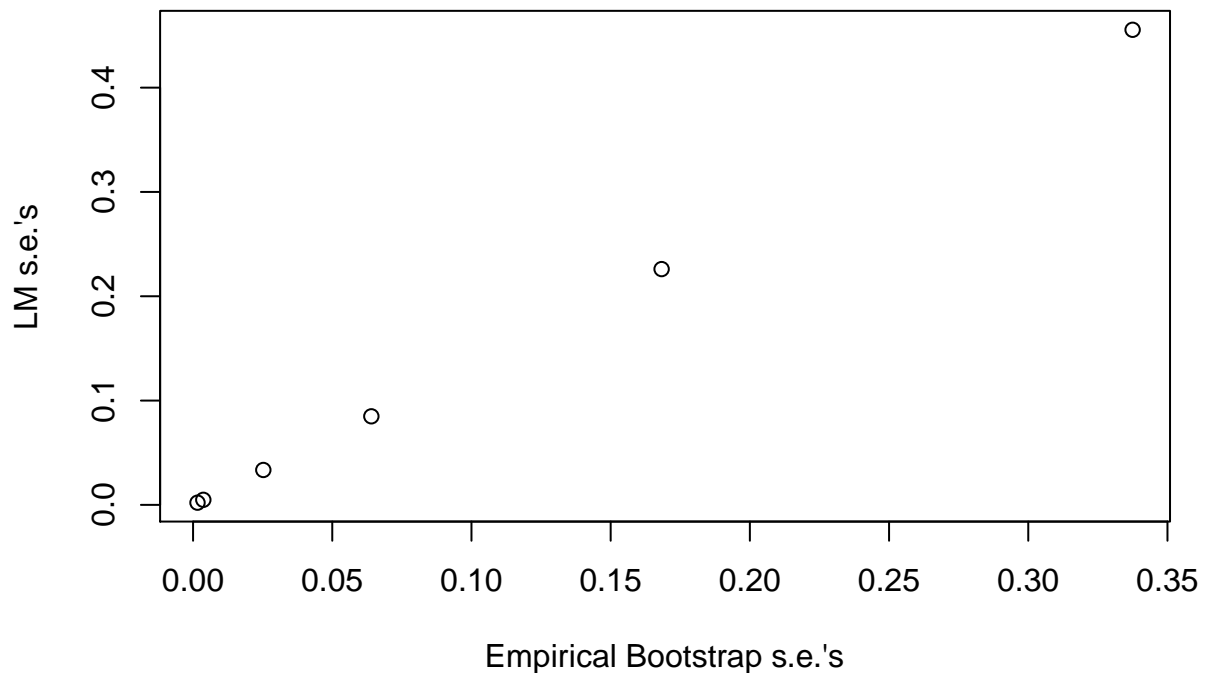
Note that the standard errors calculated in the bootstrap procedure are quite different and (lower than) the standard errors relying on likelihood theory using the `lm` function:

```
summary(fm1)$coefficients
```

```
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -3.482259e+03 8.904204e+02 -3.9108029 0.0035604037
## GNP.deflator  1.506187e-02 8.491493e-02  0.1773760 0.8631408328
## GNP          -3.581918e-02 3.349101e-02 -1.0695163 0.3126810611
## Unemployed   -2.020230e-02 4.883997e-03 -4.1364274 0.0025350917
## Armed.Forces -1.033227e-02 2.142742e-03 -4.8219853 0.0009443668
```

```
## Population  -5.110411e-02 2.260732e-01 -0.2260511 0.8262117958
## Year        1.829151e+00 4.554785e-01  4.0158898 0.0030368033
```

```
# summary(fm1)$coefficients[, "Std. Error"]
plot(y = summary(fm1)$coefficients[, "Std. Error"][-1],
     x = se_boot[-1],
     ylab= "LM s.e.'s", xlab = "Empirical Bootstrap s.e.'s"
)
```



Here,  $n$  is simply too small. We need more than 16 observations on these 7 parameters to get the bootstrap to work. With too little data, double (identical) observations create a dependency structure in the data that leads to smaller standard errors.

## 2.3 Bootstrapping Quakes and the Spearman correlation coefficient

We will use the data set “quakes” in this exercise (type in the Console `?quakes` to learn more about these data). Somehow, we are interested in the distribution of the Spearman correlation coefficient between the variables longitude (long) and latitude (lat) for 1000 seismic events. To calculate the Spearman correlation coefficient use the function `cor()` with argument `method = “spearman”`.

Draw  $B = 3000$  bootstrap samples of size  $n$  with replacement from the original data.

a)

In each of the 3000 bootstrap samples calculate the Spearman correlation coefficient and collect these coefficients in a vector (`CorSp_bs`).

**Answer:**

Option 1:

```

B <- 3e3; n <- nrow(quakes)
set.seed(20141125)
boots <- sample(1:n, n*B, replace = TRUE)
length(boots)

## [1] 3000000

Boot.lon <- matrix(quakes$long[boots], nrow = n)
Boot.lat <- matrix(quakes$lat[boots], nrow = n)

CorSp_bs <- sapply(1:B, FUN = function(i){
  cor(Boot.lon[,i], Boot.lat[,i], method = 'spearman')
})
length(CorSp_bs)

```

```
## [1] 3000
```

Option 2:

```

B <- 3e3; n <- nrow(quakes)
set.seed(20141125)
bootindex <- matrix(sample(1:n, n*B, replace = TRUE), nrow = n)
CorSp.bs2 <- sapply(1:B, FUN = function(i){
  cor(quakes$long[bootindex[,i]], quakes$lat[bootindex[,i]], method = 'spearman')
})
all.equal(CorSp_bs, CorSp.bs2)

```

```
## [1] TRUE
```

The length of the vector with bootstrapped correlations equals B, the number of bootstrap samples. The vectors created with Option 1 and 2 are identical.

b)

Make a histogram of the distribution of the vector created in **a**), and add a vertical linear to indicate the observed Spearman correlation `obs_spcor`

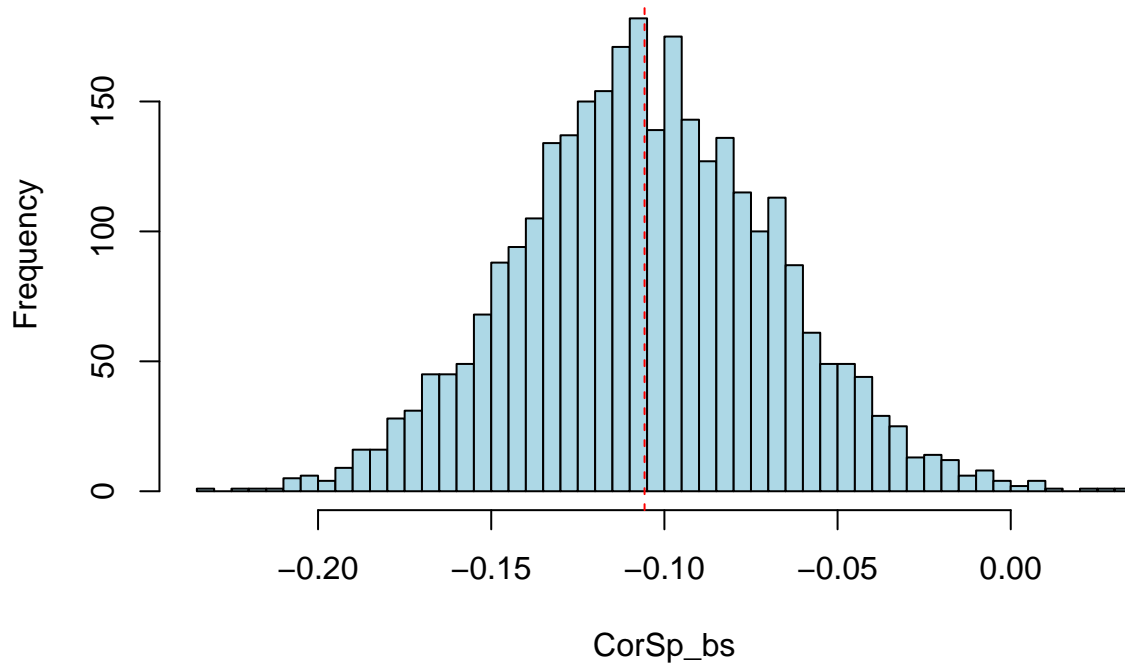
**Answer:**

```

obs_spcor <- cor(quakes$long, quakes$lat, method = 'spearman')
hist(CorSp_bs, breaks = "FD", col = 'lightblue')
abline(v = obs_spcor, lty = 2, col = 'red')

```

## Histogram of CorSp\_bs



c)

What would be the 95% confidence interval for the true Spearman correlation coefficient when using the percentiles only of the empirical bootstrapped values?

**Answer:**

```
L_perc <- quantile(CorSp_bs, 0.025)
names(L_perc) <- NULL
U_perc <- quantile(CorSp_bs, 0.975)
names(U_perc) <- NULL
c(L_perc = L_perc, U_perc = U_perc)
```

```
##      L_perc      U_perc
## -0.17670334 -0.03128189
```

d)

What is the 95% confidence interval based on the normal approximation. For the standard error of the spearman correlation coefficient, use the estimate of the standard error (= standard deviation) that you can obtain from your empirical bootstrap replicates of the spearman correlation coefficients.

**Answer:**

```
L_nrm = obs_spcor - qnorm(0.975)*sd(CorSp_bs)
U_nrm = obs_spcor - qnorm(0.025)*sd(CorSp_bs)
c(L_nrm = L_nrm, U_nrm = U_nrm)
```

```
##      L_nrm      U_nrm
## -0.17835024 -0.03310029
```

e)

Also create the 95% confidence interval based on the approximate pivot function while using the empirical bootstrap procedure (slide 59 of the lectures).

**Answer**

```
# L_piv <- obs_spcor - (quantile(CorSp_bs - obs_spcor, 0.975))
L_piv <- obs_spcor - (quantile(CorSp_bs - mean(CorSp_bs), 0.975))
names(L_piv) <- NULL
# U_piv <- obs_spcor - (quantile(CorSp_bs - obs_spcor, 0.025))
U_piv <- obs_spcor - (quantile(CorSp_bs - mean(CorSp_bs), 0.025))
names(U_piv) <- NULL
interval_piv <- c(L_piv = L_piv, U_piv = U_piv)
interval_piv
```

```
##      L_piv      U_piv
## -0.1794888 -0.0340674
```

f)

When comparing these intervals altogether, would you conclude with approx. 95% confidence that there is a (very) small negative spearman correlation between longitude and latitude regarding seismographic events?

**Answer:**

Comparing them all together:

```
c(L_perc = L_perc, U_perc = U_perc, width = U_perc - L_perc)
```

```
##      L_perc      U_perc      width
## -0.17670334 -0.03128189  0.14542145
```

```
c(L_nrm = L_nrm, U_nrm = U_nrm, width = U_nrm - L_nrm)
```

```
##      L_nrm      U_nrm      width
## -0.17835024 -0.03310029  0.14524995
```

```
c(L_piv = L_piv, U_piv = U_piv, width = U_piv - L_piv)
```

```
##      L_piv      U_piv      width
## -0.1794888 -0.0340674  0.1454214
```

Based on the confidence interval results, it seems that the conclusion may hold.

## Remaing Exercies / Self-Study (Difficult!)

### 3.1

a

According to R, the closest estimate to the ‘true’ median of our median reaction time experiment of Exercises 1.1. would be:

```
alpha <- 3; beta <- 1 # true values
med_mc_approx <- qgamma(0.5, alpha, beta)
med_mc_approx
```

```
## [1] 2.67406
```

Come up with an estimate of the coverage for the three types of 95% confidence intervals for the observed median reaction time of Exercise 1.1 for as well the parametric bootstrap and the empirical bootstrap. What is the proportion taken over e.g.  $B_{MC} = 1000$  confidence intervals for  $B_{boot} = 1e3$  bootstrap replicates that each confidence interval envelopes the true parameter for the median, i.e. `med_mc_approx`?

The three types of confidence intervals are:

1. the 95% confidence interval using only the quantiles of the paramteric bootstrap replicates:

- `L(median_observed) = quantile(t_pboots, 0.025)`
- `U(median_observed) = quantile(t_pboots, 0.975)`

2. the 95% confidence interval using the normal approximation, which will be calculated as follows (L = lower bound; U = Upper bound):

- `L(median_observed) = median_observed - qnorm(0.975)*sd(t_pboots)`
- `U(median_observed) = median_observed - qnorm(0.025)*sd(t_pboots)`

3. the 95% confidence interval using the quantiles of the parametric bootstrap replicates and the pivot function with scale equal to 1, which is calculated as follows:

- `L(median_observed) = median_observed - (quantile(t_pboots - mean(t_pboots), 0.975))`
- `U(median_observed) = median_observed - (quantile(t_pboots - mean(t_pboots), 0.025))`

**Note that for good estimates of the coverage of each confidence interval a much larger value for  $B_{MC}$  would be needed to reduce simulation error on the coverage estimate. Thus, we cannot form strong conclusions on our results! Also note that our definition of the coverage is also way to blunt (but workable).**

**Answer:**

For this specific example you would find that the confidence intervals using the quantiles of the bootstrapped sampling distribution does best.

The “true” median:

```
alpha <- 3; beta <- 1
med_mc_approx <- qgamma(0.5, alpha, beta)
n <- 1000 # sample size
```

For the empirical bootstrap (Warning the code may take a while!)

```
set.seed(20181121)
cov_norm <- cov_piv <- cov_perc <- 0
B_mc <- 1e3; b_mc <- 0

while (b_mc < B_mc) {
  b_mc <- b_mc + 1

  # just to notify R is working:
  if (b_mc %% 100 == 0) cat("b_mc = ", b_mc, "\n")

  b <- 0;
  X <- rgamma(n, alpha, beta) # data
  t_obs <- median(X)
  B <- 1e3; t_eboots <- rep(NA, B)
  while (b < B) {
    b <- b + 1
    X_b <- sample(X, replace = TRUE)
    t_eboots[b] <- median(X_b)
  }

  c1_perc <- quantile(t_eboots, 0.025) < med_mc_approx
  c2_perc <- quantile(t_eboots, 0.975) > med_mc_approx
  cov_perc <- cov_perc + as.numeric(c1_perc & c2_perc)

  c1_norm <- t_obs - qnorm(0.975) * sd(t_eboots) < med_mc_approx
  c2_norm <- t_obs - qnorm(0.025) * sd(t_eboots) > med_mc_approx
  cov_norm <- cov_norm + as.numeric(c1_norm & c2_norm)

  c1_piv <- t_obs - quantile(t_eboots - mean(t_eboots), 0.975) < med_mc_approx
  c2_piv <- t_obs - quantile(t_eboots - mean(t_eboots), 0.025) > med_mc_approx
  cov_piv <- cov_piv + as.numeric(c1_piv & c2_piv)
}

## b_mc = 100
## b_mc = 200
## b_mc = 300
## b_mc = 400
## b_mc = 500
## b_mc = 600
## b_mc = 700
## b_mc = 800
## b_mc = 900
## b_mc = 1000
```

```
cov_perc / b_mc
```

```
## [1] 0.944
```

```
cov_norm / b_mc
```

```
## [1] 0.93
```

```
cov_piv / b_mc
```

```
## [1] 0.916
```

**b**

Instead of using the empirical bootstrap, would you be able to repeat the whole experiment to estimate the coverage for each of these three interval types while using the parametric bootstrap?

**Answer:**

Just for fun, if we would like to use the parametric bootstrap the estimate of the coverage is:

```
set.seed(20181122)
cov_norm <- cov_piv <- cov_perc <- 0
B_mc <- 1e3; b_mc <- 0

while (b_mc < B_mc) {
  b_mc <- b_mc + 1

  # just to notify R is working:
  if (b_mc %% 100 == 0) cat("b_mc = ", b_mc, "\n")

  b <- 0;
  X <- rgamma(n, alpha, beta) # data
  t_obs <- median(X)
  beta_hat <- mean(X) / var(X)
  alpha_hat <- mean(X) * beta_hat
  B <- 1e3; t_pboots <- rep(NA, B)
  while (b < B) {
    b <- b + 1
    X_b <- rgamma(n, alpha_hat, beta_hat)
    t_pboots[b] <- median(X_b)
  }

  c1_piv <- t_obs - quantile(t_pboots - mean(t_pboots), 0.975) < med_mc_approx
  c2_piv <- t_obs - quantile(t_pboots - mean(t_pboots), 0.025) > med_mc_approx
  cov_piv <- cov_piv + as.numeric(c1_piv & c2_piv)

  c1_norm <- t_obs - qnorm(0.975) * sd(t_pboots) < med_mc_approx
  c2_norm <- t_obs - qnorm(0.025) * sd(t_pboots) > med_mc_approx
  cov_norm <- cov_norm + as.numeric(c1_norm & c2_norm)

  c1_perc <- quantile(t_pboots, 0.025) < med_mc_approx
  c2_perc <- quantile(t_pboots, 0.975) > med_mc_approx
  cov_perc <- cov_perc + as.numeric(c1_perc & c2_perc)
}
```



```
## b_mc = 100
## b_mc = 200
## b_mc = 300
## b_mc = 400
## b_mc = 500
## b_mc = 600
## b_mc = 700
## b_mc = 800
## b_mc = 900
## b_mc = 1000
```

```
cov_norm / b_mc
```

```
## [1] 0.942
```

```
cov_piv / b_mc
```

```
## [1] 0.941
```

```
cov_perc / b_mc
```

```
## [1] 0.986
```

### 3.2. Bootstrapping: a task from an Old Exam (SCR 2010)

a)

Consider a study with following factors:

- **sex**: factor with levels “male” and “female”
- **treat**: factor with levels “active” and “placebo”
- **age**: factor with levels “young” and “old”
- **bmi**: factor with levels “under”, “normal”, “over” and “obese”

Construct the `data.frame` that contains all possible combinations of the factors. *Hint: the dimension are 32 rows and 4 columns*

**Answer:**

```
# all possible combinations of the factors in a data.frame
dat <- data.frame(expand.grid(
  sex = factor(0:1, labels=c("male", "female")),
  treat = factor(0:1, labels= c("active", "placebo")),
  age = factor(0:1, labels=c("young", "old")),
  bmi = factor(0:3, labels=c("under", "normal", "over", "obese"))
))
```

b)

For each possible combination of factor levels in (a) simulate 50 objects from the linear regression model

$$y_i \sim N(\mu_i, \sigma^2), \quad i = 1, \dots, n = 50 \times 32,$$

where  $\sigma = 3$  and

$$\mu_i = \beta_0 + \beta_1 \text{female}_i + \beta_2 \text{placebo}_i + \beta_3 \text{old}_i + \beta_4 \text{normal}_i + \beta_5 \text{over}_i + \beta_6 \text{obese}_i$$

with  $\text{female}_i$  denoting the dummy variable for females,  $\beta_2 \text{placebo}_i$  the dummy variable for placebo patients,  $\beta_3 \text{old}_i$  the dummy variable for old patients,  $\beta_4 \text{normal}_i$  the dummy variable for patients with normal BMI,  $\beta_5 \text{over}_i$  the dummy variable for obese patients.

For the regression coefficients take the values  $\beta_0 = 10, \beta_1 = 1, \beta_2 = -2, \beta_3 = -1.3, \beta_4 = 0.1, \beta_5 = -1, \beta_6 = -1.2$ .

*Hint: check the function `model.matrix` and its use in predicting and fitting linear models*

**Answer:**

```
sigma <- 2
betas <- c(b0 = 10, b1 = -1, b2 = -2, b3 = -1.3,
           b4 = .1, b5 = -1, b6 = -1.2)
X <- model.matrix(~dat$sex + dat$treat + dat$age + dat$bmi)
mu.i <- X %*% betas
set.seed(123)
BMI <- dat[rep(rownames(dat), each = 50), ]
BMI$y <- rnorm(50*32)*sigma + rep(mu.i, each = 50)
```

c)

Using the data you simulated in (b), fit the following four linear regression models:

- M1 additive model: include only the main effects of `sex`, `treat`, `age`, and `bmi`.
- M2 2-way interaction model: include only the main effects of `sex`, `treat`, `age`, and `bmi`, and all the 2-way interaction terms.
- M3 3-way interaction model: include only the main effects of `sex`, `treat`, `age`, and `bmi`, all the 2-way interaction terms and all the 3-way interaction terms.
- M4 4-way interaction model: include only the main effects of `sex`, `treat`, `age`, and `bmi`, all the 2-way interaction terms, all the 3-way interaction terms, and all the 4-way interaction terms. The formula  $y \sim (x_1 + x_2 + x_3 + x_4)^4$  might be helpful.

Also, extract the design matrices for each model. Remember from linear regression:  $\mathbf{b} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$ . The  $\mathbf{X}$  is the design matrix.

*Stronger hint: really, check the function `model.matrix()`*

**Answer:**

```
M1 <- lm(y ~ sex + treat + age + bmi, dat = BMI)
M2 <- lm(y ~ (sex + treat + age + bmi)^2, dat = BMI)
M3 <- lm(y ~ (sex + treat + age + bmi)^3, dat = BMI)
```

```
M4 <- lm(y ~ (sex + treat + age + bmi)^4, dat = BMI)
# extract design matrices
X.mod1 <- model.matrix(M1)
X.mod2 <- model.matrix(M2)
X.mod3 <- model.matrix(M3)
X.mod4 <- model.matrix(M4)
```

OR, when tired of writing down all 4 models...

```
for (i in 2:4) {
  mname <- paste0("M", i)
  assign(mname, lm(
    formula = bquote(y ~ (sex + treat + age + bmi)^.(i)),
    dat = BMI
  ))
  modmat <- paste0("X.mod", i)
  assign(modmat, model.matrix(get(mname)))
}
```

d)

We are interested in quantifying the predictive ability of these models. As a measure of predictive ability we will use the adjusted  $R^2$  defined as

$$R_{adj}^2 = 1 - (1 - R^2) \frac{n - 1}{n - p - 1}$$

where  $n$  denotes the sample size,  $p$  the total number of regressors in the linear model (but not counting the constant term / intercept), and

$$R^2 = COR(\hat{y}, y)^2,$$

the squared correlation between the predicted outcome  $\hat{y}$  and the outcome variable  $y$ .

Use the Bootstrap to estimate a validated  $R^2$  for each of the four models using 200 samples with replacement from the original sample. Is your validated  $R^2$  closer to the observed  $R^2$  or the observed  $R_{adj}^2$ ? Do you see any differences when the number of parameter  $P$  becomes larger?

It is often claimed in biostatistics and in the social sciences that  $R_{adj}^2$  is more valid to use for predictive ability than  $R^2$ . So, one would expect that the validated values for  $R^2$  are closer to those of the observed  $R_{adj}^2$  when the number of parameter  $P$  becomes larger. Is this the case?

**Answer:**

```
# function to compute (adjusted) R^2 based on the design
# matrix, the response vector and the estimated
# coefficients
R2 <- function(X, y, betas){
  f <- c(X %*% betas)
  R2 <- cor(f, y)^2
  return(R2)
}
```

```
adjR2 <- function(X, y, betas){
  R2 <- R2(X,y, betas)
  n <- nrow(X)
  cnst <- (n - 1)/(n - length(betas))
  out <- 1 - (1 - R2)*cnst
  return(out)
}
```

```
set.seed(321)
y <- BMI$y
n <- length(y)
B <- 100 # how many Bootstrap samples
# to save time, define bootstrapped indices beforehand:
ind <- matrix(sample(1:n, n*B, replace = TRUE), n, B)

# to store the bootstrapped R2's:
r2.boot.mods <- matrix(numeric(B*4), nrow = B)

# starting the procedure:
for (i in 1:B) {
  # index of the i-th Bootstrap sample
  ind.i <- ind[, i]

  # using the following code in the loop takes too much time!!
  # ind.i <- sample(1:n, n*B, replace = TRUE)

  # fit the models again
  mod1.new <- lm.fit(X.mod1[ind.i, ], y[ind.i])
  mod2.new <- lm.fit(X.mod2[ind.i, ], y[ind.i])
  mod3.new <- lm.fit(X.mod3[ind.i, ], y[ind.i])
  mod4.new <- lm.fit(X.mod4[ind.i, ], y[ind.i])

  # compute R2
  r2.boot.mods[i, 1] <- R2(X.mod1, y, mod1.new$coefficients)
  r2.boot.mods[i, 2] <- R2(X.mod2, y, mod2.new$coefficients)
  r2.boot.mods[i, 3] <- R2(X.mod3, y, mod3.new$coefficients)
  r2.boot.mods[i, 4] <- R2(X.mod4, y, mod4.new$coefficients)
}
colMeans(r2.boot.mods) # all validated R2 values
```

```
## [1] 0.3464615 0.3479096 0.3462501 0.3474711
```

Resulting in the following changes

```
colMeans(r2.boot.mods) - c(
  R2(X.mod1, y, M1$coefficients),
  R2(X.mod2, y, M2$coefficients),
  R2(X.mod3, y, M3$coefficients),
  R2(X.mod4, y, M4$coefficients)
)
```

```
## [1] -0.001789259 -0.006726841 -0.010638202 -0.011821913
```

Thus, the validated  $R^2$  values are indeed lower...

```
colMeans(r2.boot.mods) - c(
  adjR2(X.mod1, y, M1$coefficients),
  adjR2(X.mod2, y, M2$coefficients),
  adjR2(X.mod3, y, M3$coefficients),
  adjR2(X.mod4, y, M4$coefficients)
)
```

```
## [1] 0.0006655403 0.0006207511 0.0008240063 0.0008451256
```

and closer to the adjusted  $R^2$  values, especially when  $P$  becomes larger.