

## Exercise 6

1. Consider the following linear regression model

$$Y_i = \beta_1 + \beta_2 x_{i2} + \beta_3 x_{i3} + \epsilon_i, \quad i = 1, \dots, 100, \quad \beta_1 = 1, \beta_2 = -2, \beta_3 = 3, \quad (1)$$

where the pairs  $x_{i2}, x_{i3}$  lie on a  $\{1, \dots, 10\} \times \{1, \dots, 10\}$ -grid, i.e.,

```
x2 <- rep(1:10, 10)
x3 <- rep(1:10, each = 10)
```

In this exercise, we will simulate datasets from this model using different error distributions and perform linear regression. Confidence intervals for the regression parameters will then be estimated using classical “normal theory” and bootstrapping.

- a) Implement your own bootstrap routine as an R function that takes a data frame with columns `x2`, `x3` and `y` as input and that returns three confidence intervals, one for each regression parameter  $\beta_j$ ,  $j = 1, 2, 3$ .
- b) Simulate 100 datasets<sup>1</sup> from model (1) computing each time classical “normal theory” 0.95-confidence intervals and bootstrap<sup>2</sup> 0.95-confidence intervals for the three regression parameters.

For the simulations, use three different types of error distributions (resulting in 300 simulated datasets, all in all):

1.  $\epsilon_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1)$ ,
2.  $\epsilon_i \stackrel{\text{i.i.d.}}{\sim} t_3$  (R function `rt`),
3.  $\epsilon_i = e_i - 1$ , where  $e_i \stackrel{\text{i.i.d.}}{\sim} \text{Exp}(1)$  (R function `rexp`).

How often do the confidence intervals include the true values (estimated coverage rate)?

**R-hints:** To make your results reproducible, use `set.seed(11)` at the beginning of your simulation experiment.

Classical confidence intervals for output objects of `lm` can be computed using `confint`. To sample the bootstrap-indices for your own bootstrap-routine, use the functions `sample` and/or `replicate` (Look at the help-files!).

Use  $B = 1000$  bootstrap samples. Beware that bootstrapping is computationally quite expensive. In order to avoid long waiting times, first develop and test your code with few simulations and bootstrap samples (say, 10 each), and augment these numbers only when your code works.

---

<sup>1</sup>It depends on the computer time you can spend whether you try 50, 100, 200 or 1000 simulations. It may need lots of time, because each time a complete bootstrap simulation has to be carried out. You can always downsize your simulations by simulating fewer datasets and/or varying the number of bootstrap replicates.

<sup>2</sup>The bootstrap replicates should be generated by sampling from the set  $\{(x_1, Y_1), \dots, (x_{100}, Y_{100})\}$ , and **not** by resampling the residuals as in the model-based bootstrap.

- c) Repeat the bootstrap calculation of the confidence intervals by using the function `boot.ci` from package `boot`.

**R-hints:** The function `boot` from package `boot` allows automatic bootstrapping of statistics on given data. To apply this function, you have to write your own R-function which returns the regression coefficients and has arguments `dat` and `ind`. `dat` is a data frame containing the variables `y`, `x2` and `x3`, and `ind` is a vector of indices (see help page, parameter `statistic`). Such a function may look like this:

```
lmcoefs <- function(dat, ind) coef(lm(y~x2+x3, data=dat[ind,]))
```

Then use the `boot` function:

```
bst.sample <- boot(data=dat, statistic=lmcoefs, R=B)
```

Bootstrap confidence intervals are then computed by `boot.ci`, which may look as follows:

```
bst.ci <- boot.ci(bst.sample, conf=1-alpha, type="basic", index=k)
```

`bst.sample` is the output of `boot`, `index` should be 1 for the intercept parameter, 2 and 3 for the regression parameters (if computed as in `lmcoefs` above). The interval bounds come as values `bst.ci$basic[4]` and `bst.ci$basic[5]`.

- d) In this part of the exercise we want to compare the usual  $L_1$ -loss  $\frac{1}{n} \sum_{i=1}^n |y_i - \hat{m}(x_i)|$  with the  $L_1$ -generalization error  $\mathbf{E}[|Y_{\text{new}} - \hat{m}(X_{\text{new}})|]$ . This time the  $L_1$ -generalization-error is estimated by bootstrapping instead of cross-validation as described in the manuscript. Do 100 simulations for each of the given error distributions. In each simulation calculate the two quantities of interest and compare their averages over the whole range of simulations. A histogram of the two quantities may be informative too. You might want to recycle the bootstrap-samples you generated above.

**Preliminary discussion:** Friday, April 23, 2010.

**Deadline:** Friday, April 30, 2010.