

## Model Evaluation

The task of this exercise is to evaluate predictive performance of several linear models and choose the best one for predicting medical costs billed by the health insurance. The prediction influences monthly payments of the clients. Therefore, a high predictive accuracy is a priority. The file `insurance-all.Rdata` contains information on the covariates and dependent variable which are described in detail in Table 1. The simulated data set was created using demographic statistics from the U.S. Census Bureau for the covariates, and thus approximately reflect real-world conditions.<sup>1</sup>

Table 1: Description of the Variables

Variable	Description
age	age of primary beneficiary (numeric)
sex	insurance contractor gender (factor: “female”, “male”)
bmi	body mass index, ideally 18.5 to 24.9 (numeric)
children	number of children covered by health insurance / number of dependents (numeric)
smoker	smoking (factor: “yes”, “no”)
region	the beneficiary’s residential area in the US (factor: “northeast”, “southeast”, “southwest”, “northwest”)
charges	individual medical costs billed by health insurance

### Group Home Assignment (max. 4 points)

The mandatory group home assignment has to be submitted before 12:00 o’clock prior PC-session 2. It is obligatory to solve the assignment in R Markdown with `echo = TRUE` option for every code chunk and generate a PDF file with the solution. Generating an HTML file from the R Markdown and converting it into the PDF is also possible. Make sure that the final PDF file has no readability issues. The PDF with the answers to the six questions below as well as the file with the R Markdown code has to be submitted via Canvas.

Download the data set `insurance-all.Rdata` from Canvas. Load the data into R. Install and load the package `ggplot`.

1. How many observations are in the whole dataset? How many covariates were collected? (0.5 points)
2. What is the highest number of children who are covered by one health insurance? (0.5 points)
3. Look at the percentage shares of smokers and non-smokers within each region, e.g. in one region there could be 80% smokers and 20% non-smokers. Which region has the lowest share of smokers? What is the share of smokers there? (0.5 points)

<sup>1</sup>Lantz, B. (2013) Machine Learning with R. Packt Publishing.

4. Create a scatter plot with **charges** on the  $y$ -axis and **age** on the  $x$ -axis. Distinguish in the plot by color which data points belong to smokers and non-smokers (coded in covariate **smoker**). Describe the patterns in the data. (0.5 points)
5. Write a function that has a data argument and three string arguments **x.variable**, **y.variable**, **color.variable** which are used to generate a scatter plot of two variables against each other and use the third to color the points. Use this function to plot **x.variable="bmi"**, **y.variable="charges"**, **color="sex"** and interpret the results. (1 point)  
*Hint: the aesthetic mapping `aes()` in `ggplot` does not interpret string inputs. Search for an alternative.*
6. Write a function that creates a boxplot of the **bmi** variable split by another variable that is passed as the argument. Do you see any difference in the **bmi** by region? (1 point)