

Data Science

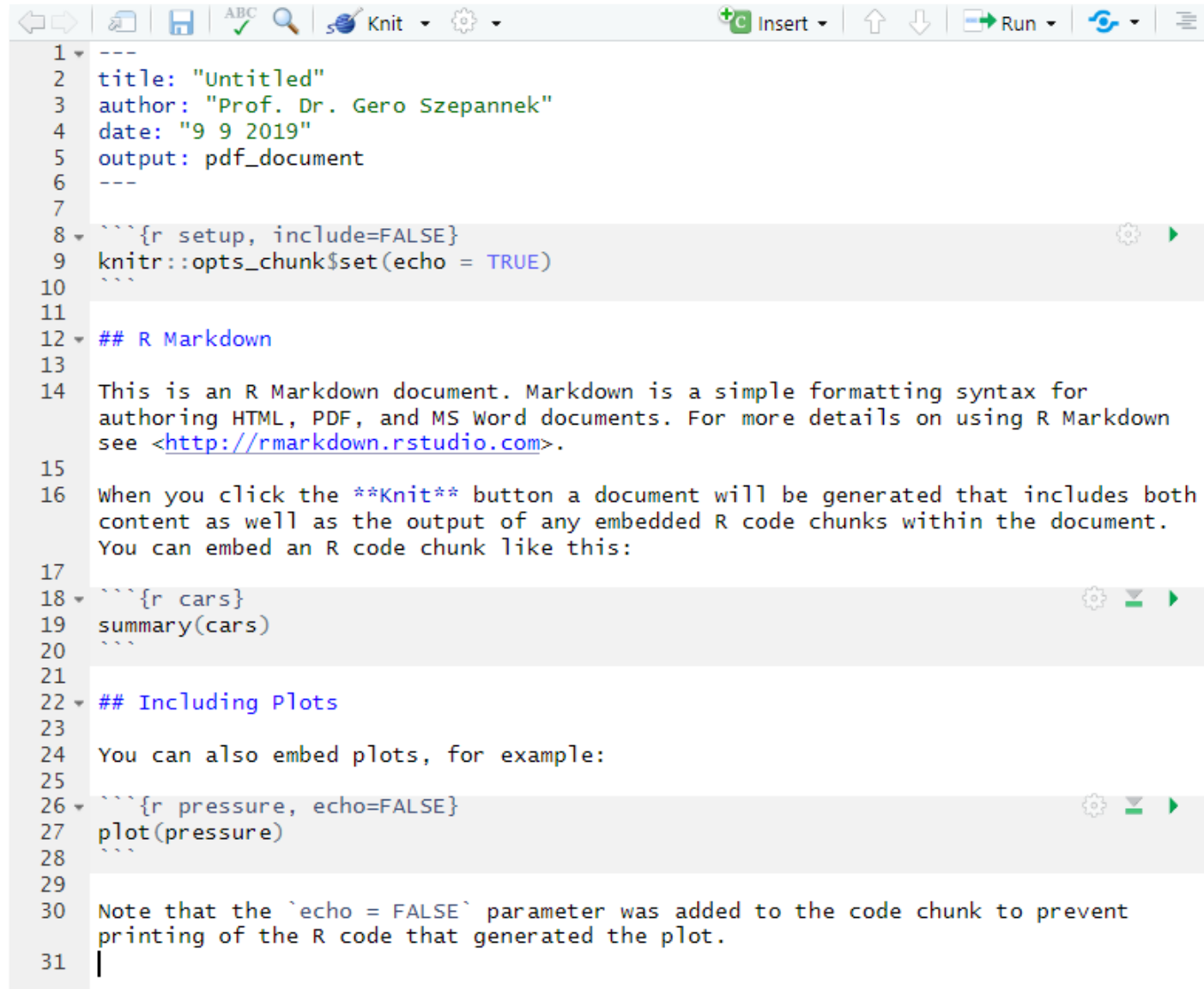
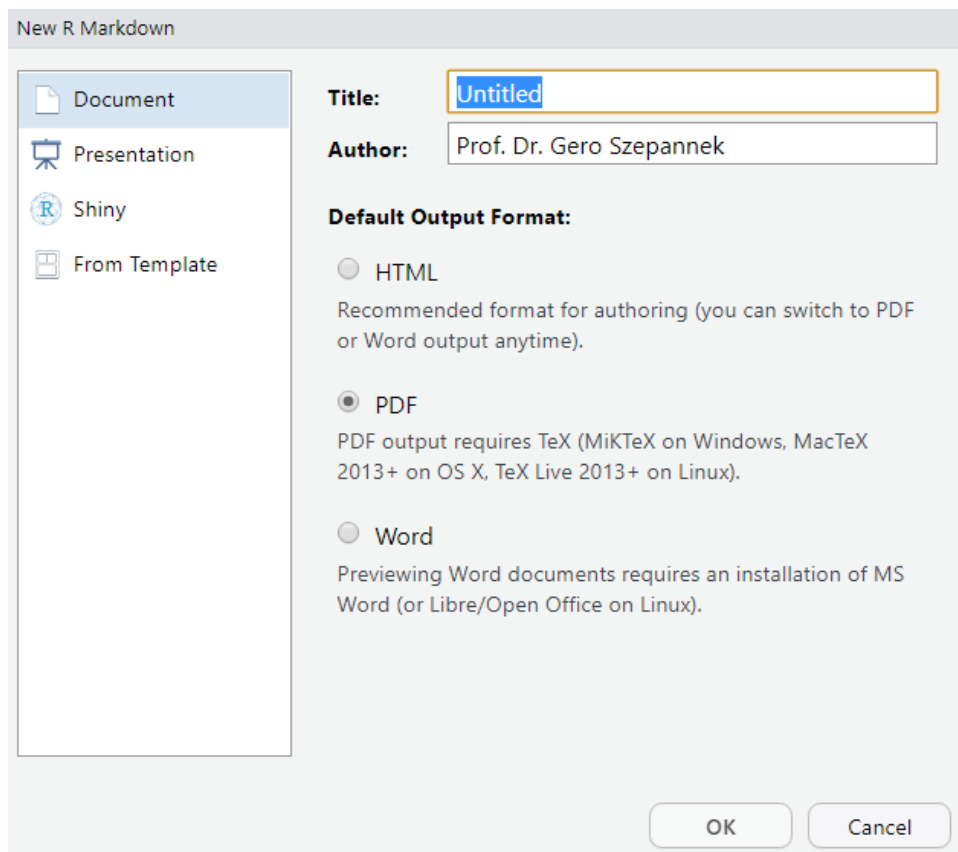
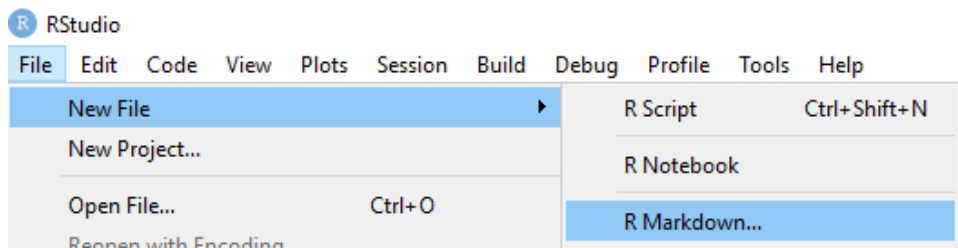
Turning Analyses directly into Results

Gero Szepannek



Literate Programming

“Let us change our traditional attitude to the construction of programs: Instead of imagining that our main task is to instruct a computer what to do, let us concentrate rather on explaining to human beings what we want a computer to do.” (D. Knuth (1984): Literate Programming, The Computer Journal 27 p. 97.)



```
1 ---
2 title: "Untitled"
3 author: "Prof. Dr. Gero Szepannek"
4 date: "9 9 2019"
5 output: pdf_document
6 ---
7
8 ```{r setup, include=FALSE}
9 knitr::opts_chunk$set(echo = TRUE)
10 ```
11
12 ## R Markdown (Sub-)Section header
13
14 This is an R Markdown document. Markdown is a simple formatting syntax for
15 authoring HTML, PDF, and MS Word documents. For more details on using R Markdown
16 see <http://rmarkdown.rstudio.com>.
17
18 When you click the Knit button a document will be generated that includes both
19 content as well as the output of any embedded R code chunks within the document.
20 You can embed an R code chunk like this:
21
22 ```{r cars}
23 summary(cars)
24 ```
25
26 ## Including Plots
27
28 You can also embed plots, for example:
29
30 ```{r pressure, echo=FALSE}
31 plot(pressure)
32 ```
33
34 Note that the `echo = FALSE` parameter was added to the code chunk to prevent
35 printing of the R code that generated the plot.
```

Meta-information

Link

R chunk

R chunk where code is not shown
in output but only the figure

bold
italic

<https://rmarkdown.rstudio.com/>

<https://bookdown.org/yihui/rmarkdown/>

Exercise

1. Create a Rmarkdown template as shown on the previous slides!
2. Run the code find out what documents are created!
3. Modify the template such that a boxplot of the variable age of the custdata.csv is displayed.

Exercise / Case Study: Identify Customers who don't have a health insurance

Assignment of teams:

1. Import the file custdata.csv!
2. Create a Rmarkdown file to present your results!
3. How many observations have the data?
4. What variables are in the data? Which variable is the target variable that indicates whether a customer has a health insurance?
5. What percentage of customers has no health insurance?
6. What can you tell about missing data (NA)? (Note: the function `is.na()` checks for missings.)
7. What can you tell about the variables income and age with regard to plausibility?
8. Create a table to compare whether the percentages of insured persons varies by the state (state.of.res). Do also consider absolute numbers. What would you conclude?
9. Create some boxplots and mosaicplots: which variable can be used to identify persons that don't have a health insurance?
10. **Submit your team's results as a .Rmd that should run on my computer.**