

Applied Data Analysis, Übung 4

HENRY HAUSTEIN

Vorbereitung

```
1  install.packages("dplyr")
2  library(dplyr)
3  install.packages("readxl")
4  library(readxl)
5  data = read_excel("data.xlsx", na = "NA")
6
7  data$gender = factor(data$gender)
8  levels(data$gender) = c("male", "female", "diverse")
9  data$employment = factor(data$employment)
10 levels(data$employment) = c("student", "employed", "unemployed")
11 data$education = factor(data$education)
12 levels(data$education) = c("no degree", "secondary", "intermediate",
13                             "high school", "academic")
13 data$play_frequency = factor(data$play_frequency)
14 levels(data$play_frequency) = c("never", "every few months", "
    every few weeks", "1-2 days a week", "3-5 days a week", "daily
    ")
15 data$treatment = factor(data$treatment)
16 levels(data$treatment) = c("control", "lootbox in task reward", "
    lootbox picture", "badge")
17 data$age = sapply(data$age, function(year) {2016-year})
18 data$rt6 = as.numeric(data$rt6)
19 data$rt7 = as.numeric(data$rt7)
20 data$rt8 = as.numeric(data$rt8)
21 data$rt9 = as.numeric(data$rt9)
22 data$rt10 = as.numeric(data$rt10)
23 data$rt11 = as.numeric(data$rt11)
24 data$rt12 = as.numeric(data$rt12)
25 data$rt13 = as.numeric(data$rt13)
26 data$rt14 = as.numeric(data$rt14)
27 subControl = subset(data, treatment == "control")
```

Task 1

```
1  install.packages("ggplot2")
2  library(ggplot2)
3  ggplot(data, aes(x = age, y = 'total time'))
```

```

4 + geom_point()
5 + geom_smooth(method = "lm")
6 lm_tt_age = lm('total time' ~ age, data = data)
7 summary(lm_tt_age)
8 ggplot(subControl, aes(x = age, y = 'total time'))
9 + geom_point()
10 + geom_smooth(method = "lm")
11 lm_tt_age_control = lm('total time' ~ age, data = subControl)
12 summary(lm_tt_age_control)
13
14 ggplot(data, aes(x = tasks_completed, y = 'total time'))
15 + geom_point()
16 + geom_smooth(method = "lm")
17 lm_tt_tc = lm('total time' ~ tasks_completed, data = data)
18 summary(lm_tt_tc)
19 ggplot(subControl, aes(x = tasks_completed, y = 'total time'))
20 + geom_point()
21 + geom_smooth(method = "lm")
22 lm_tt_tc_control = lm('total time' ~ tasks_completed, data =
    subControl)
23 summary(lm_tt_tc_control)
24
25 ggplot(data, aes(x = rt0, y = 'total time'))
26 + geom_point()
27 + geom_smooth(method = "lm")
28 lm_tt_rt0 = lm('total time' ~ rt0, data = data)
29 summary(lm_tt_rt0)
30 ggplot(subControl, aes(x = rt0, y = 'total time'))
31 + geom_point()
32 + geom_smooth(method = "lm")
33 lm_tt_rt0_control = lm('total time' ~ rt0, data = subControl)
34 summary(lm_tt_rt0_control)

```

Das führt zu folgenden Modellen:

	Full Model	Control Group
(Intercept)	230665.09*** (53287.10)	298612.42 (156127.98)
age	3635.06 (2032.53)	-1457.79 (5666.45)
R ²	0.01	0.00
Adj. R ²	0.01	-0.01
Num. obs.	437	99

*** $p < 0.001$; ** $p < 0.01$; * $p < 0.05$

	Full Model	Control Group
(Intercept)	170842.92*** (44899.68)	280409.31* (128200.11)
tasks_completed	14859.12*** (4095.99)	-2518.17 (15008.84)
R ²	0.03	0.00
Adj. R ²	0.03	-0.01
Num. obs.	437	99

*** $p < 0.001$; ** $p < 0.01$; * $p < 0.05$

	Full Model	Control Group
(Intercept)	157695.17*** (5821.01)	119595.89*** (5304.52)
rt0	1.01*** (0.01)	0.99*** (0.01)
R ²	0.91	0.99
Adj. R ²	0.91	0.99
Num. obs.	437	99

*** $p < 0.001$; ** $p < 0.01$; * $p < 0.05$

Man sieht, dass offensichtlich nur die Variable `rt0` die Variable `total time` gut erklären kann.

Task 2

```

1 lm_tt_age_rt0 = lm('total time' ~ age + rt0, data = data)
2 summary(lm_tt_age_rt0)
3
4 lm_tt_age_rt0_interaction = lm('total time' ~ age + rt0 + age:rt0,
5                               data = data)
6 summary(lm_tt_age_rt0_interaction)
7
8 lm_tt_treatment_rt0 = lm('total time' ~ treatment + treatment:rt0
9                           + 0, data = data)
10 summary(lm_tt_treatment_rt0)

```

Das führt zu folgenden Modellen:

	without interaction terms	with interaction terms
(Intercept)	107465.81*** (15600.65)	115139.46*** (20861.24)
age	2048.80*** (591.55)	1700.27* (863.35)
rt0	1.01*** (0.01)	0.98*** (0.06)
age:rt0		0.00 (0.00)
R ²	0.92	0.92
Adj. R ²	0.92	0.92
Num. obs.	437	437

*** $p < 0.001$; ** $p < 0.01$; * $p < 0.05$

	Model 1
treatmentcontrol	119595.89*** (10859.57)
treatmentlootbox in task reward	224714.96*** (12861.18)
treatmentlootbox picture	135344.63*** (18221.65)
treatmentbadge	153816.42*** (10165.92)
treatmentcontrol:rt0	0.99*** (0.02)
treatmentlootbox in task reward:rt0	1.03*** (0.04)
treatmentlootbox picture:rt0	1.06*** (0.13)
treatmentbadge:rt0	1.01*** (0.02)
R ²	0.96
Adj. R ²	0.96
Num. obs.	437

*** $p < 0.001$; ** $p < 0.01$; * $p < 0.05$

Sobald wir Interaktionsterme erlauben, ist das Alter der Probanden gar nicht mehr so wichtig. Die Behandlung scheint aber einen großen Einfluss auf die `total time` zu haben.