

Applied Data Analysis, Übung 5

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)
```

Task 1

```
1  aov_tc_treatment = aov(tasks_completed ~ treatment, data = data)
2  summary(aov_tc_treatment)
3
```

```

4 pairwise.t.test(data$tasks_completed, data$treatment, p.adjust.
  method = "none")
5 pairwise.t.test(data$tasks_completed, data$treatment, p.adjust.
  method = "bonferroni")
6 pairwise.t.test(data$tasks_completed, data$treatment, p.adjust.
  method = "holm")

```

Ergebnis von ANOVA ist

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
treatment	3	1585.61	528.54	34.24	$< 2 \cdot 10^{-16}$
Residuals	433	6683.72	15.44		

Da ANOVA auf

- $H_0 : \mu_1 = \dots = \mu_g$
- $\exists i, j : \mu_i \neq \mu_j$

testet, sehen wir am p-value, dass die Treatments einen unterschiedlichen Mittelwert haben. Die paarweisen t-Tests haben folgendes Ergebnis (ohne p-value-adjustment, [Bonferroni-Korrektur](#)¹, [Bonferroni-Holm Korrektur](#)²):

	control	lootbox in task reward	lootbox picture
lootbox in task reward	$< 2 \cdot 10^{-16}$, $< 2 \cdot 10^{-16}$, $< 2 \cdot 10^{-16}$		
lootbox picture	0.05967, 0.35800, 0.05967	$7.3 \cdot 10^{-14}$, $4.4 \cdot 10^{-13}$, $3.6 \cdot 10^{-13}$	
badge	$2.8 \cdot 10^{-8}$, $1.7 \cdot 10^{-7}$, $1.1 \cdot 10^{-7}$	$1.9 \cdot 10^{-5}$, 0.00011, $5.6 \cdot 10^{-5}$	0.00013, 0.00079, 0.00026

Wir sehen, dass bei fast allen Gruppen sich die Mittelwerte unterscheiden, bis auf $\mu_{\text{lootbox pic}} = \mu_{\text{control}}$.

Task 2

```

1 install.packages("car")
2 library(car)
3 data %>% group_by(treatment) %>% summarise(completeTasks_var = var
  (tasks_completed))
4 leveneTest(data$tasks_completed, data$treatment)
5 oneway.test(tasks_completed ~ treatment, data = data, var.equal =
  FALSE)
6
7 qqnorm(data$tasks_completed)
8 qqline(data$tasks_completed)
9 kruskal.test(tasks_completed ~ treatment, data = data)

```

Die Varianzen unterscheiden sich:

¹ $p^* = \frac{\alpha}{n}$ mit n Anzahl der Tests

² $p^* = \frac{\alpha}{m+1-k}$ mit p_1, \dots, p_m geordneten p-values und für die k -te These

	treatment	completeTasks_var
1	control	12.43
2	lootbox in task reward	12.29
3	lootbox picture	15.56
4	badge	19.82

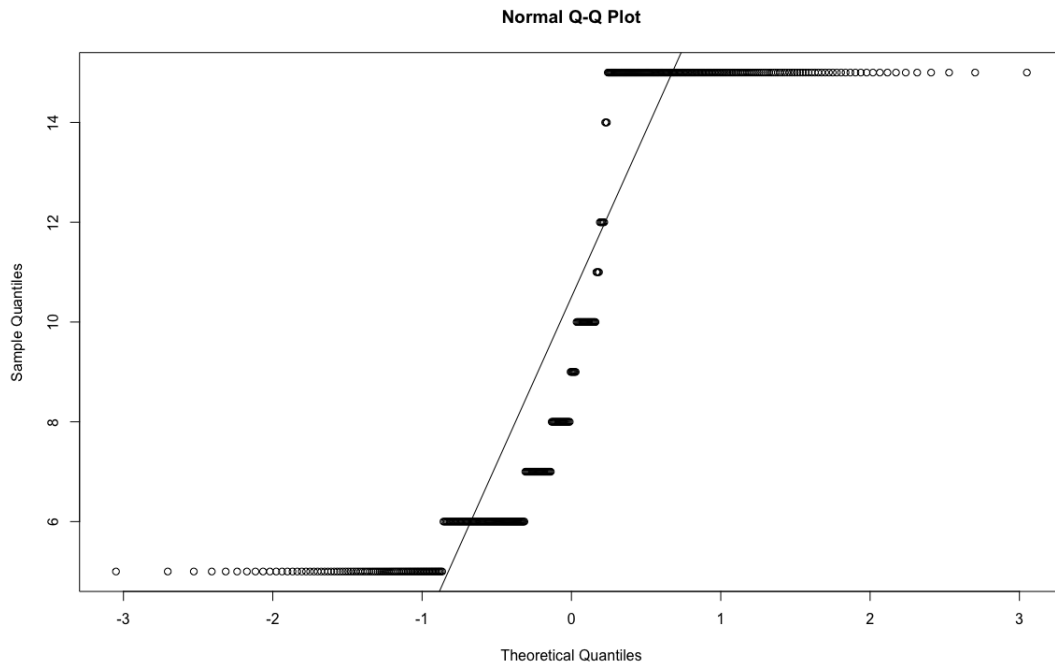
Der Levene Test testet

- $H_0 : \sigma_1^2 = \dots = \sigma_g^2$
- $H_1 : \exists i, j : \sigma_i^2 \neq \sigma_j^2$

und hat folgendes Ergebnis

	Df	F value	Pr(>F)
group	3	14.14	$< 2 \cdot 10^{-16}$
	433		

Ja, die Varianzen unterscheiden sich. Der Oneway-ANOVA-Test hat einen p-value von $< 2 \cdot 10^{-16}$, also gibt es auch hier Unterschiede zwischen den Gruppen.



Offensichtlich kann man hier auch nicht von Normalverteilung ausgehen. Der Kruskal-Wallis-Test testet, ob die Gruppen aus der gleichen Population (H_0) kommen. Das scheint hier nicht der Fall zu sein, der p-value ist $1.51 \cdot 10^{-15}$.