

Applied Data Analysis, Übung 3

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
28 subsetControl = data %>% filter(treatment == "control")
29 subsetLootPic = data %>% filter(treatment == "lootbox picture")
30 subsetLootInTask = data %>% filter(treatment == "lootbox in task
    reward")
31 subsetBadge = data %>% filter(treatment == "badge")
32 subsetYoung = data %>% filter(old == FALSE)
33 subsetOld = data %>% filter(old == TRUE)
```

Task 1

```
1 t.test(data$'total time', mu = 320000, alternative = "two.sided",
  conf.level = 0.95)
2 t.test(subsetYoung$'total time', mu = 320000, alternative = "two.
  sided", conf.level = 0.95)      t.test(subsetOld$'total time',
  mu = 320000, alternative = "two.sided", conf.level = 0.95)
3
4 t.test(subsetYoung$tasks_completed, mu = mean(subsetOld$tasks_
  completed), alternative = "two.sided", conf.level = 0.95)
5
6 pairwise.t.test(data$tasks_completed, data$treatment, p.adjust.
  method = "none", pool.sd = FALSE)
```

Bei allen Personen und den jungen Personen ist der p-value unter 5%, das heißt wir können H_0 nicht ablehnen. Bei den alten Personen ist der p-value bei 50%, wir lehnen H_0 ab.

Bei allen andern Tests sind die p-values unter 5%, auch da lehnen wir H_0 ab. Die Mittelwerte sind also verschieden voneinander.

Task 2

```
1 install.packages("ggplot2")
2 library(ggplot2)
3
4 ggplot(data, aes(x = 'total time'))
5 + geom_histogram(aes(y = ..density..))
6 + stat_function(fun = dnorm, args = list(mean = mean(data$'total
  time'), sd = sd(data$'total time'), color = "red")
7 ggplot(data, aes(x = log('total time')))
8 + geom_histogram(aes(y = ..density..))
9 + stat_function(fun = dnorm, args = list(mean = mean(log(data$'
  total time')), sd = sd(log(data$'total time'))), color = "red"
  )
10 ggplot(data, aes(sample = 'total time'))
11 + geom_qq()
12 + geom_qq_line()
13 ggplot(data, aes(sample = log('total time')))
14 + geom_qq()
15 + geom_qq_line()
16
17 ggplot(data, aes(x = tasks_completed))
18 + geom_histogram(aes(y = ..density..))
19 + stat_function(fun = dnorm, args = list(mean = mean(data$tasks_
  completed), sd = sd(data$tasks_completed)), color = "red")
20 ggplot(data, aes(x = log(tasks_completed)))
21 + geom_histogram(aes(y = ..density..))
22 + stat_function(fun = dnorm, args = list(mean = mean(log(data$
  tasks_completed)), sd = sd(log(data$tasks_completed))), color
  = "red")
23 ggplot(data, aes(sample = tasks_completed))
```

```

24 + geom_qq()
25 + geom_qq_line()
26 ggplot(data, aes(sample = log(tasks_completed)))
27 + geom_qq()
28 + geom_qq_line()

```

Man sieht, dass die Daten auf keinen Fall normalverteilt sind, höchstens die Variable `total time` ist etwas log-normalverteilt.

Es gibt einen Test auf Normalverteilung, den Shapiro-Wilk-Test. Nullhypothese dabei ist, dass die Daten normalverteilt sind.

```

1 shapiro.test(data$`total time`)
2 shapiro.test(data$tasks_completed)
3 shapiro.test(log(data$`total time`))
4 shapiro.test(log(data$tasks_completed))

```

Liefert uns p-values von $< 2.2 \cdot 10^{-16}$ (bzw. $4.767 \cdot 10^{-15}$ beim Logarithmus von `total time`), die Nullhypothese wird also angelehnt. Normalverteilt sind die Daten nicht.

Task 3

Der Kern des Skriptes ist der folgende Ablauf, der 10 000 mal wiederholt wird:

1. Wähle aus dem gegebenen Datensatz 90 Beobachtungen aus.
2. Bilde den Mittelwert von `tasks_completed` aus den ausgewählten Beobachtungen.
3. Speichere den Mittelwert ab.

Das Central Limit Theorem (Zentraler Grenzwertsatz ↗ *Statistik 2*) sagt nun, dass diese 10 000 Mittelwerte normalverteilt sind, was dann in dem Skript mit einigen Abbildungen gezeigt wird.

Der Prozess des mehrfachen Auswählens aus vorhandenen Daten heißt *Bootstrapping* und ist mit Vorsicht zu genießen. Man erweckt damit den Eindruck als hätte man 10 000 Personen untersucht, aber in Wahrheit kommen die Daten nur von 99 Personen (für das Bootstrapping der Kontrollgruppe) bzw. 114 Personen (für das Bootstrapping der Lootbox-Bild-Gruppe). Wenn also die Datengrundlage zu klein ist, kann man mittels Bootstrapping es so aussehen lassen, als sei die Datengrundlage riesig.

Anschließend wird versucht einen t-Test durchzuführen, indem man die 10 000 Mittelwerte aus der Kontrollgruppe mit den 10 000 Mittelwerten aus der Lootbox-Bild-Gruppe vergleicht. Das 95%-Konfidenzintervall für den Unterschied wird mittels Quantilen bestimmt und ergibt $[-2.1, 0.1]$. Führt man den richtigen t-Test durch, so erhält man:

```

1 t.test(subControl$tasks_completed, subLootPic$tasks_completed,
        paired = FALSE)

```

als Konfidenzintervall $[-2.0282492, -0.0100283]$.

Task 4

```

1 ggplot(data, aes(x = age, y = `total time`))
2 + geom_point()
3 + geom_smooth(method = "lm")

```

```

4 summary(lm('total time' ~ age, data = data))
5 ggplot(subsetControl, aes(x = age, y = 'total time'))
6 + geom_point()
7 + geom_smooth(method = "lm")
8 summary(lm('total time' ~ age, data = subsetControl))
9
10 ggplot(data, aes(x = tasks_completed, y = 'total time'))
11 + geom_point()
12 + geom_smooth(method = "lm")
13 summary(lm('total time' ~ tasks_completed, data = data))
14 ggplot(subsetControl, aes(x = tasks_completed, y = 'total time'))
15 + geom_point()
16 + geom_smooth(method = "lm")
17 summary(lm('total time' ~ tasks_completed, data = subsetControl))
18
19 ggplot(data, aes(x = rt0, y = 'total time'))
20 + geom_point()
21 + geom_smooth(method = "lm")
22 summary(lm('total time' ~ rt0, data = data))
23 ggplot(subsetControl, aes(x = rt0, y = 'total time'))
24 + geom_point()
25 + geom_smooth(method = "lm")
26 summary(lm('total time' ~ rt0, data = subsetControl))

```

Der Residual Standard Error ist schwer von der Größe her zu interpretieren, aber der R^2 sollte nahe bei 1 liegen. Nur für das Modell `total time ~ rt0` ist das so, dort liegt er bei 0.9139 bzw. 0.9906 für die Kontrollgruppe. Für die anderen beiden linearen Modelle verbessert sich durch die Einschränkung auf die Kontrollgruppe nichts.