

# Übungszettel Logistische Regression, Exercise Sheet Logistic Regression

M.Psy.205, Dozent: Peter Zezula

Johannes Brachem ([johannes.brachem@stud.uni-goettingen.de](mailto:johannes.brachem@stud.uni-goettingen.de))

30 Mai, 2022 21:39

## Deutsche Version

### Links

[Übungszettel als PDF-Datei zum Drucken](#)

[Übungszettel mit Lösungen](#)

[Lösungszettel als PDF-Datei zum Drucken](#)

[Der gesamte Übungszettel als .Rmd-Datei \(Zum Downloaden: Rechtsklick > Speichern unter...\)](#)

### Hinweise zur Bearbeitung

1. Bitte beantworten Sie die Fragen in einer .Rmd Datei. Sie können Sie über [Datei > Neue Datei > R Markdown...](#) eine neue R Markdown Datei erstellen. Den Text unter dem *Setup Chunk* (ab Zeile 11) können Sie löschen. [Unter diesem Link](#) können Sie auch unsere Vorlage-Datei herunterladen (Rechtsklick > Speichern unter...).
2. Informationen, die Sie für die Bearbeitung benötigen, finden Sie auf der [Website der Veranstaltung](#)
3. Zögern Sie nicht, im Internet nach Lösungen zu suchen. Das effektive Suchen nach Lösungen für R-Probleme im Internet ist tatsächlich eine sehr nützliche Fähigkeit, auch Profis arbeiten auf diese Weise. Die beste Anlaufstelle dafür ist der [R-Bereich der Programmiererplattform Stackoverflow](#)
4. Auf der Website von R Studio finden Sie sehr hilfreiche Übersichtszettel zu vielen verschiedenen R-bezogenen Themen. Ein guter Anfang ist der [Base R Cheat Sheet](#)

### Ressourcen

Da es sich um eine praktische Übung handelt, können wir Ihnen nicht alle neuen Befehle einzeln vorstellen. Stattdessen finden Sie hier Verweise auf sinnvolle Ressourcen, in denen Sie für die Bearbeitung unserer Aufgaben nachschlagen können.

Ressource	Beschreibung
Field, Kapitel 8	Buchkapitel, das Schritt für Schritt erklärt, worum es geht, und wie man logistische Regressionen in R durchführt. <b>Große Empfehlung!</b>

## Tipp der Woche

Wenn Sie Zeilenumbrüche in Ihrem Code verwenden, ist er für Sie und Andere deutlich besser lesbar. Sie können z.B. nach jeder Pipe (%>%) oder nach dem + in ggplot2-Befehlen sinnvolle Zeilenumbrüche einbauen. Zögern Sie auch nicht, längere einzelne Befehle in mehrere Zeilen umzubrechen: Nach jedem Komma kann eine neue Zeile beginnen, wenn Sie es möchten.

Es kann sein, dass die Einrückung bei mehrzeiligen Befehlen irgendwann unordentlich wird. Wenn das passiert, markieren Sie den Code und drücken strg + i (windows), bzw. cmd + i (mac). Ihr Code wird dadurch automatisch eingerückt.

## 1 Daten einlesen

1. Laden Sie das tidyverse und fügen Sie die entsprechende Code-Zeile an den Anfang Ihrer .Rmd-Datei ein.
2. Setzen Sie ein sinnvolles Arbeitsverzeichnis und fügen Sie die entsprechende Code-Zeile an den Anfang Ihrer .Rmd-Datei ein.
3. Laden Sie den Datensatz stats\_lectures unter [diesem Link](#) herunter und speichern ihn in Ihrem Arbeitsverzeichnis, dort z.B. in einem Unterordner data.
4. Lesen Sie den Datensatz in R ein.

### Erklärung der Variablen

Jede Zeile enthält Daten von einem/r StudentIn. Beobachtet wurde, ob die Studierenden an Statistik-Lehrveranstaltungen teilnehmen. Erhoben wurden dazu die Temperatur, wie sehr die Studis Sonne genießen und wie sehr sie Statistik genießen. Letzteres wurde mit einem Fragebogen gemessen, der aus 5 Items besteht.

**Hinweis:** Die Daten sind vollständig fiktiv und spiegeln natürlich auch nicht die Meinung der Lehrkräfte wieder. Dem Tutor ist nur gerade kein besseres Beispiel eingefallen.

Name	Bedeutung
attend_or_not	Hat zwei Ausprägungen: "course not attended" und "course attended". Die Ausprägungen geben an, ob der/die Studi zur Lehrveranstaltung gegangen ist.
temperature	Die Temperatur. Höher ist heißer.
sun_joy	Wie sehr genießen die Studis Sonne. Höher ist stärkerer Genuss
stats_joy1 bis stats_joy5	Wie sehr genießen die Studis Statistik. Höher ist stärkerer Genuss.

### Lösung

```
library(tidyverse)
```

### Unteraufgabe 1

**Unteraufgabe 2** Beachten Sie hier, dass in R alle Ordner mit / getrennt werden. Deshalb können Sie Windows-Dateipfade nicht einfach kopieren, diese sind mit \ getrennt. Ersetzen Sie dort \ durch /.

Der Code ist:

```
setwd()
```

Hier ein Beispiel mit meinem Dateipfad.

```
# Beispiel  
setwd("~/ownCloud/_Arbeit/Hiwi Peter/gitlab_sheets/public")
```

**Unteraufgabe 3** Bitte folgen Sie den Anweisungen in der Aufgabenstellung.

```
#stats_lectures <- read_csv("public/mv/data/stats_lectures.csv")  
stats_lectures <- read_csv("https://md.psych.bio.uni-goettingen.de/mv/data/div/stats_lectures.csv")
```

**Unteraufgabe 4**

```
## Rows: 400 Columns: 8  
## -- Column specification -----  
## Delimiter: ","  
## chr (1): attend_or_not  
## dbl (7): temperature, sun_joy, stats_joy1, stats_joy2, stats_joy3, stats_joy4, stats_joy5  
##  
## i Use `spec()` to retrieve the full column specification for this data.  
## i Specify the column types or set `show_col_types = FALSE` to quiet this message.
```

## 2 Datenaufbereitung

1. Bilden Sie für jede Person im Datensatz den Mittelwert aus den 5 Items stats\_joy1 bis stats\_joy5 und fügen Sie diesen als neue Variable dem Datensatz hinzu. *Hinweis: Eine Google-Suche kann hier sehr hilfreich sein. Achten Sie auch auf Kleinigkeiten in den Befehlen, die Sie finden.* Falls Sie diese Aufgabe nicht lösen können, aber dennoch weiter machen möchten, finden Sie hier funktionierende Syntax, die Sie in Ihr Script kopieren können. [Link](#)
2. Nutzen Sie den Befehl `factor()`, um die Variable `attend_or_not` in einen Faktor umzuformen. Verwenden Sie das Argument `levels = c()`, um die Faktorstufen anzugeben. So stellen Sie sicher, dass Sie wissen, welche Kategorie die Baseline darstellt. Falls Sie diese Aufgabe nicht lösen können, aber dennoch weiter machen möchten, finden Sie hier funktionierende Syntax, die Sie in Ihr Script kopieren können. [Link](#)

## Lösung

**Unteraufgabe 1** Der Befehl `mean()` muss hier mit `mean(c())` verwendet werden, damit die Zeilenweisen Mittelwerte richtig gebildet werden. Wir verwenden `mutate()`, um eine neue Variable zu erstellen.

```
stats_lectures <- stats_lectures %>%  
  rowwise() %>% # wir wollen zeilenweise vorgehen  
  mutate(stats_joy = mean(c(stats_joy1, stats_joy2, # mittelwert bilden  
                           stats_joy3, stats_joy4,  
                           stats_joy5)))
```

```
stats_lectures <- stats_lectures %>% mutate(attend_or_not = factor(attend_or_not, levels = c("course no
```

## Unteraufgabe 2

### 3 Erste logistische Regression

1. Nutzen Sie den Befehl `glm()` mit dem Argument `family = binomial()`, um die Kursbesuche der Studierenden mit deren Genuss von Statistik vorherzusagen. Benutzen Sie dafür den Mittelwert, den Sie in der vorherigen Aufgabe gebildet haben, als Prädiktor.
2. Lassen Sie sich den Output mit `summary()` anzeigen. "Null deviance" gibt die Deviance-Statistik für das Null-Modell an, "Residual Deviance" für das Alternativmodell (user spezifiziertes). Was können Sie aus der Deviance ableiten? (Siehe Field, Kapitel 8.3.2: Assessing the model: the deviance statistic)
3. Bilden Sie mit `exp()` das Exponential zur Basis e für die Regressionskoeffizienten. Dies ist das Odds-Ratio. *Tipp: Mit \$coefficients können Sie direkt auf die Koeffizienten zugreifen, wenn Sie den Namen, den Sie dem Modell gegeben haben, vor das \$-Zeichen schreiben.*
4. Interpretieren Sie das Odds Ratio für den vorliegenden Fall. (Siehe Field, Kapitel 8.3.6: The odds ratio)
5. Wenden Sie `confint()` auf das Modell an, um die Konfidenzintervalle für die Prädiktoren zu erhalten. Bilden Sie auch hiervon das Exponential mit `exp()`, um die Konfidenzintervalle für die Odds-Ratios zu erhalten.
6. Erstellen Sie einen Plot, um Ihr Modell zu visualisieren.
  - a) Wenden Sie den Befehl `fitted()` auf das Modell an, um die durch das Modell vorhergesagten Wahrscheinlichkeiten für Ihre Studierenden zu erhalten.
  - b) Fügen Sie diese Werte Ihrem Datensatz hinzu.
  - c) Nutzen Sie `ggplot`, um einen Plot zu erstellen. Verwenden Sie auf der X-Achse `stats_joy` und auf der y-Achse die in a) und b) erstellten vorhergesagten Wahrscheinlichkeiten.

## Lösung

```
model1 <- glm(attend_or_not ~ stats_joy,  
               data = stats_lectures, family = binomial())
```

## Unteraufgabe 1

```
summary(model1)
```

## Unteraufgabe 2

```
##  
## Call:  
## glm(formula = attend_or_not ~ stats_joy, family = binomial(),  
##       data = stats_lectures)  
##  
## Deviance Residuals:
```

```

##      Min       1Q   Median       3Q      Max
## -1.7081   0.7275   0.7275   0.7275   0.7275
##
## Coefficients: (1 not defined because of singularities)
##                 Estimate Std. Error z value Pr(>|z|)
## (Intercept)    1.1942    0.1184   10.09 <2e-16 ***
## stats_joy      NA        NA        NA        NA
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 433.82 on 399 degrees of freedom
## Residual deviance: 433.82 on 399 degrees of freedom
## AIC: 435.82
##
## Number of Fisher Scoring iterations: 4

```

Die Deviance ist ein Indikator dafür, wie viel Information nicht durch das Modell aufgeklärt werden kann. Eine größere Deviance bedeutet, dass weniger Information erklärt werden kann. Mehr dazu können Sie hier nachlesen: Field, Kapitel 8.3.2: Assessing the model: the deviance statistic.

```
exp(model1$coefficients)
```

### Unteraufgabe 3

```

## (Intercept)  stats_joy
##     3.301075          NA

```

**Unteraufgabe 4** Wenn hier bei Studi A die Freude an Statistik um einen Punkt höher ist als bei Studi B, dann heißt das, dass die Chance, dass Studi A die Statistik-Vorlesung besucht, ca. 2.86-mal größer ist, als dass Studi B es tut.

D.h. eine Steigerung der Freude an Statistik um einen Punkt führt zu einer 2.86-mal höheren Chance, die Statistik-Vorlesung zu besuchen.

```
exp(confint(model1))
```

### Unteraufgabe 5

```

## Waiting for profiling to be done...

##           2.5 %   97.5 %
## (Intercept) 2.629478 4.184281
## stats_joy      NA        NA

```

## Unteraufgabe 6

- a) Ich nutze hier `head()`, um nur die ersten 6 Werte anzeigen zu lassen. Andernfalls würde das Dokument sehr voll werden.

```
head(fitted(model1))
```

```
##      1      2      3      4      5      6
## 0.7675 0.7675 0.7675 0.7675 0.7675 0.7675
```

b)

```
stats_lectures$pred_prob <- fitted(model1)
```

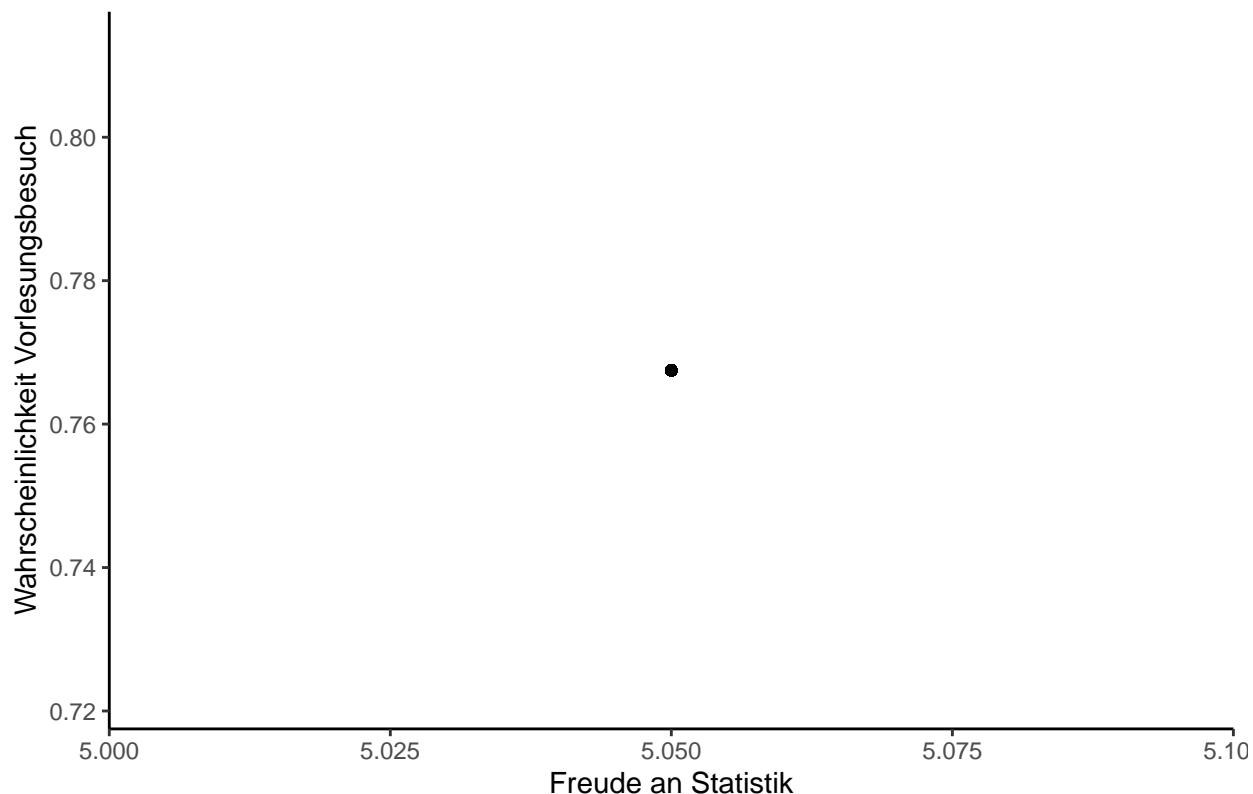
c)

```
# Objekt erstellen
plot1 <- ggplot(stats_lectures, aes(x = stats_joy, y = pred_prob))

# Plot befüllen
plot1 <- plot1 + geom_point() +
  labs(x = "Freude an Statistik",
       y = "Wahrscheinlichkeit Vorlesungsbesuch",
       title = "Vorlesungsbesuche in Abhängigkeit von Freude an Statistik") +
  theme_classic()

# plot anzeigen
plot1
```

## Vorlesungsbesuche in Abhängigkeit von Freude an Statistik



## 4 Modellvergleich logistische Regression

1. Erstellen Sie drei weitere Modelle mit dem `glm()`-Befehl:
  - a) Fügen Sie den ersten Modell den Prädiktor `temperature` hinzu
  - b) Fügen Sie dem Modell aus a) den Prädiktor `sun_joy` hinzu
  - c) Fügen Sie dem Modell aus b) die Interaktion aus `temperature` und `sun_joy` als Prädiktor hinzu.
2. Testen Sie mit dem Befehl `anova()`, welches der Modelle am besten auf die Daten passt. Nutzen Sie die Option `test = "Chisq"`, um einen p-Wert zu erhalten.
  - a) Sehen Sie sich mit `summary()` den Output für das beste Modell an.
  - b) Welche Hypothesen würden Sie aus dieser Analyse ableiten?
3. Berechnen Sie für das beste Modell die Odds Ratios der Prädiktoren und deren Konfidenzintervalle.
4. Berechnen Sie für das beste Modell die drei Pseudo-R<sup>2</sup> Indices (Siehe Field, Kapitel 8.3.3: Assessing the model: R und R<sup>2</sup>). Sie können dafür die unten stehende Funktion nutzen. *Mehr dazu in Field, R's Souls' Tip 8.2.*
  - a) Kopieren Sie den Code und führen Sie ihn aus.
  - b) Sie können nun die Funktion `logisticPseudoR2s()` auf ihr Modell anwenden.

```
logisticPseudoR2s = function(LogModel) {  
  dev = LogModel$deviance  
  nullDev = LogModel$null.deviance  
  modelN = length(LogModel$fitted.values)  
  R.1 = 1 - dev / nullDev
```

```

R.cs = 1 - exp(-(nullDev - dev) / modelN)
R.n = R.cs / (1 - (exp(-(nullDev / modelN)))))

cat("Pseudo R^2 for logistic regression\n")
cat("Hosmer and Lemeshow R^2 ", round(R.l, 3), "\n")
cat("Cox and Snell R^2 ", round(R.cs, 3), "\n")
cat("Nagelkerke R^2 ", round(R.n, 3), "\n")
}

```

## Lösung

```

# das erste Modell
model1 <- glm(attend_or_not ~ stats_joy,
               data = stats_lectures, family = binomial())

# a)
model2 <- glm(attend_or_not ~ stats_joy + temperature,
               data = stats_lectures, family = binomial())

# b)
model3 <- glm(attend_or_not ~ stats_joy + temperature + sun_joy,
               data = stats_lectures, family = binomial())

# c) (kann auf zwei Arten geschrieben werden)
model4 <- glm(attend_or_not ~ stats_joy + temperature*sun_joy,
               data = stats_lectures, family = binomial())

model4 <- glm(attend_or_not ~ stats_joy + temperature +
               sun_joy + temperature:sun_joy,
               data = stats_lectures, family = binomial())

```

## Unteraufgabe 1

```
anova(model1, model2, model3, model4, test = "Chisq")
```

## Unteraufgabe 2

```

## Analysis of Deviance Table
##
## Model 1: attend_or_not ~ stats_joy
## Model 2: attend_or_not ~ stats_joy + temperature
## Model 3: attend_or_not ~ stats_joy + temperature + sun_joy
## Model 4: attend_or_not ~ stats_joy + temperature + sun_joy + temperature:sun_joy
##   Resid. Df Resid. Dev Df Deviance Pr(>Chi)
## 1       399    433.82
## 2       398    433.20  1     0.622    0.4303
## 3       397    433.12  1     0.083    0.7727

```

```

## 4      396    385.71  1   47.405 5.774e-12 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

Modell 4 passt mit deutlichem Abstand am besten auf die Daten.

a)

```
summary(model4)
```

```

##
## Call:
## glm(formula = attend_or_not ~ stats_joy + temperature + sun_joy +
##       temperature:sun_joy, family = binomial(), data = stats_lectures)
##
## Deviance Residuals:
##    Min      1Q  Median      3Q     Max
## -2.3668  0.2048  0.6041  0.6903  1.7613
##
## Coefficients: (1 not defined because of singularities)
##                  Estimate Std. Error z value Pr(>|z|)
## (Intercept)      1.43536   0.13866 10.352 < 2e-16 ***
## stats_joy        NA         NA         NA         NA
## temperature     0.05659   0.14162  0.400   0.689
## sun_joy        -0.01291   0.14656 -0.088   0.930
## temperature:sun_joy -0.99906   0.16895 -5.913 3.35e-09 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 433.82 on 399 degrees of freedom
## Residual deviance: 385.71 on 396 degrees of freedom
## AIC: 393.71
##
## Number of Fisher Scoring iterations: 5

```

b) Man kann folgende Hypothesen ableiten:

1. Je höher die Freude an Statistik ist, desto höher ist die Wahrscheinlichkeit, dass Studierende Statistik-Vorlesungen besuchen.
2. Es gibt eine Interaktion zwischen Temperatur und Freude an der Sonne. Möglicherweise ist es so, dass eine hohe Temperatur bei SonnenliebhaberInnen dazu führt, dass die Statistik-Veranstaltungen weniger häufig besucht werden.

```

# Odds Ratios
exp(model4$coefficients)

```

### Unteraufgabe 3

```

##          (Intercept)      stats_joy   temperature      sun_joy temperature:sun_joy
##        4.2011631           NA         1.0582247       0.9871727           0.3682260

# Konfidenzintervalle. exp() nicht vergessen!
exp(confint(model4))

## Waiting for profiling to be done...

##          2.5 %    97.5 %
## (Intercept) 3.2281763 5.5649819
## stats_joy     NA         NA
## temperature  0.8018372 1.3995432
## sun_joy       0.7398945 1.3166720
## temperature:sun_joy 0.2599905 0.5048007

```

## Unteraufgabe 4

- a) Siehe Aufgabenstellung.
- b)

```
logisticPseudoR2s(model4)
```

```

## Pseudo R^2 for logistic regression
## Hosmer and Lemeshow R^2  0.111
## Cox and Snell R^2  0.113
## Nagelkerke R^2  0.171

```

## 5 Rendern (knit)

Lassen Sie die Datei mit **Strg + Shift + K** (Windows) oder **Cmd + Shift + K** (Mac) rendern. Sie sollten nun im “Viewer” unten rechts eine “schön aufpolierte” Version ihrer Datei sehen. Falls das klappt: Herzlichen Glückwunsch! Ihr Code kann vollständig ohne Fehlermeldung gerendert werden. Falls nicht: Nur mut, das wird schon noch! Gehen Sie auf Fehlersuche! Ansonsten schaffen wir es ja in der Übung vielleicht gemeinsam.

## Literatur

*Anmerkung:* Diese Übungszettel basieren zum Teil auf Aufgaben aus dem Lehrbuch *Discovering Statistics Using R* (Field, Miles & Field, 2012). Sie wurden für den Zweck dieser Übung modifiziert, und der verwendete R-Code wurde aktualisiert.

Field, A., Miles, J., & Field, Z. (2012). *Discovering Statistics Using R*. London: SAGE Publications Ltd.

## English Version

### Links

[Exercise sheet as PDF](#)

**Exercise sheet with solutions included**

[Exercise sheet with solutions included as PD](#)

The source code of this sheet as .Rmd (Right click and “save as” to download ...)

## Some hints

1. Please try to solve this sheet in an .Rmd file. You can create one from scratch using **File > New file > R Markdown....** You can delete the text beneath *Setup Chunk* (starting from line 11). Alternatively, you can download our template file under [this link](#) (right click > save as...).
2. You'll find a lot of the important information on the [website of this course](#)
3. Please don't hesitate to search the web for help with this sheet. In fact, being able to effectively search the web for problem solutions is a very useful skill, even R pros work this way all the time! The best starting point for this is the [R section on the programming site Stackoverflow](#)
4. On the R Studio website, you'll find highly helpful [cheat sheets](#) for many of R topics. The [base R cheat sheet](#) might be a good starting point.

## Ressources

Since this is a hands-on seminar, we won't be able to present each and every new command to you explicitly. Instead, you'll find here references to helpful ressources that you can use for completing this sheets.

Ressource	Description
Field, chapter 8	Book chapter explaining step by step the why and how of logistic regression in R. <b>Highly recommended!</b>

## Hint of the week

You should use new lines in your code to make it better readable for yourself and for others. You could f. e. continue in a new line after every pipe (%>%) or after a + in a sequence of ggplot2 commands. Don't hesitate to divide a complex and long command into several lines. After every comma (,) you can continue in a new line, if you like.

If your code looks scattered because of unregular indentation, you may mark a section of code and press **strg + i** (Windows), or **cmd + i** (Mac) to automatically indent your source code.

## 1 Read data

1. Load package **tidyverse** and insert the appropriate code at the beginning of your .Rmd document.
2. Change to an adequate working directory. Take care: The render process assumes, that your working directory is set to the location of your .Rmd document.
3. Load the data **stats\_lectures** from [this link](#) and store it in your working directory or in a subfolder there. You might want to call it **data**.
4. Read the data in R. Alternatively read the data directly from the URL above.

## Explanation of the variables

Each line has the data of one student. You find, whether the student took part in a statistic lesson, the current temperature, to which extend the person enjoys sun and, on the other hand, statistics. This was measured using a 5 items questionnaire.

**Hint:** These are completely fictive data. They don't reflect the docents opinion. This is just an spontaneous idea of a student helper.

name	meaning
attend_or_not	dichotomous variable: “course not attended” and “course attended”. The values code, whether the student attendet the course or not.
temperature	high values code high temperatures
sun_joy	the higher the value the higher the joy
stats_joy1 to stats_joy5	the higher the values the more the student enjoys statistics

## Solution

```
library(tidyverse)
# or
require(tidyverse)
```

### Subtask 1

**Subtask 2** Take care: R uses / to specify pathes. So you cannot simply copy and paste Windows pathes into R. You have to replace all \ with /. Again: Remember the special behavior of the rendering with regard to working directories.

The code would be:

```
setwd()
```

Here an example with a local path. Adapt that to your needs.

```
# example
setwd("P:/R/mv")
```

**Subtask 3** Please follow the instructions in subtask 3.

```
# stats_lectures <- read_csv("data/stats_lectures.csv") # relative path, working directory has to be set
# alternatively read it from the URL
stats_lectures <- read_csv("https://md.psych.bio.uni-goettingen.de/mv/data/div/stats_lectures.csv")
```

### Subtask 4

```
## Rows: 400 Columns: 8
## -- Column specification --
## Delimiter: ","
## chr (1): attend_or_not
## dbl (7): temperature, sun_joy, stats_joy1, stats_joy2, stats_joy3, stats_joy4, stats_joy5
##
## i Use `spec()` to retrieve the full column specification for this data.
## i Specify the column types or set `show_col_types = FALSE` to quiet this message.
```

## 2 Data preparation

1. Calculate the mean of items stats\_joy1 to stats\_joy5 and store it in the dataframe resp. tibble. *Hint: A google search might help. Take also care of the details of the commands, you find.* If you aren't able to solve this part of the task and want to continue, you find a working syntax for this part under [https://md.psych.bio.uni-goettingen.de/mv/sheets/data/06\\_sheet\\_task2-1.R](https://md.psych.bio.uni-goettingen.de/mv/sheets/data/06_sheet_task2-1.R) that you might integrate in your script.
2. Use the command `factor()` to transform variable `attend_or_not` into a factor. Use the argument `levels = c()` to specify factor levels. This is how you set the correct reference group or baseline. If you want to keep on and don't know how to do this, you find a working syntax under [https://md.psych.bio.uni-goettingen.de/mv/sheets/data/06\\_sheet\\_task2-2.R](https://md.psych.bio.uni-goettingen.de/mv/sheets/data/06_sheet_task2-2.R)

### Solution

**Subtask 1** The command `mean()` has to be done like `mean(c())` to have the means calculated rowwise. We use `mutate()` to create a new variable.

```
stats_lectures <- stats_lectures %>%
  rowwise() %>% # wir want rowwise calculation
  mutate(stats_joy = mean(c(stats_joy1, stats_joy2, # calculate means
                            stats_joy3, stats_joy4,
                            stats_joy5)))
# we could alternatively use base systems rowMeans()
# {}
```

```
stats_lectures <- stats_lectures %>% mutate(attend_or_not = factor(attend_or_not, levels = c("course no
```

### Subtask 2

## 3 First logistic regression

1. Use the command `glm()` with argument `family = binomial()` to predict, which students attend the statistics course and whether they enjoy it. The predictor is the mean, that you calculated recently.
2. Check the output using `summary()`. “Null deviance” is the deviance statistic for the null model, “residual deviance” is the same for the user defined alternative model. What can we learn from the deviance? C. f. Field (2012), chapter 8.3.2: Assessing the model: the deviance statistic.
3. Use `exp()` to calculate the exponent to base e for the regression coefficient. The result is the odds ratio.  
*Tip: \$coefficients gives you direct access to the coefficients. As always you have to put the name of your result model and a \$ before.*
4. Interpret the odds ratio for your example. (See Field (2012), chapter 8.3.6: The odds ratio)
5. Apply `confint()` to your result model to get confidence intervals for your predictors. Also use `exp()` to get the confidence intervals for the odds ratios.
6. Make a plot to visualize your model.
  - a) Apply `fitted()` to your model to get the predicted probabilities for your student subjects.
  - b) Add these values to your data set.
  - c) Use `ggplot()` to make a plot. Put `stats_joy` to the x-axes and on the y-axes the predicted probabilities, you calculated in a) and b).

## Solution

```
model1 <- glm(attend_or_not ~ stats_joy,
               data = stats_lectures, family = binomial())
```

### Subtask 1

```
summary(model1)
```

### Subtask 2

```
##  
## Call:  
## glm(formula = attend_or_not ~ stats_joy, family = binomial(),  
##       data = stats_lectures)  
##  
## Deviance Residuals:  
##      Min        1Q    Median        3Q       Max  
## -1.7081    0.7275    0.7275    0.7275    0.7275  
##  
## Coefficients: (1 not defined because of singularities)  
##                 Estimate Std. Error z value Pr(>|z|)  
## (Intercept)  1.1942     0.1184   10.09  <2e-16 ***  
## stats_joy      NA         NA         NA         NA  
## ---  
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1  
##  
## (Dispersion parameter for binomial family taken to be 1)  
##  
## Null deviance: 433.82 on 399 degrees of freedom  
## Residual deviance: 433.82 on 399 degrees of freedom  
## AIC: 435.82  
##  
## Number of Fisher Scoring iterations: 4
```

The deviance is an indicator of how much information can *not* be explained by our model. The bigger the deviance, the less information can be explained. Read more about this in Field (2012), chapter 8.3.2 Assessing the model: The deviance statistic.

```
exp(model1$coefficients)
```

### Subtask 3

```
## (Intercept)  stats_joy  
##      3.301075          NA
```

**Subtask 4** If student A has a stats\_joy of 1 point higher than student B, this means, that the chance, that student A visits the statistics class is 2.86 times higher than for student B.

In other words: An increase of joy in statistics of 1 point results in a 2.86 times higher probability to visit the statistics class.

```
exp(confint(model1))
```

### Subtask 5

```
## Waiting for profiling to be done...

##           2.5 %    97.5 %
## (Intercept) 2.629478 4.184281
## stats_joy      NA        NA
```

### Subtask 6

a) We use `head()`, to see the first 6 values only. Thus we avoid an extremely large document.

```
head(fitted(model1))
```

```
##      1      2      3      4      5      6
## 0.7675 0.7675 0.7675 0.7675 0.7675 0.7675
```

b)

```
stats_lectures$pred_prob <- fitted(model1)
```

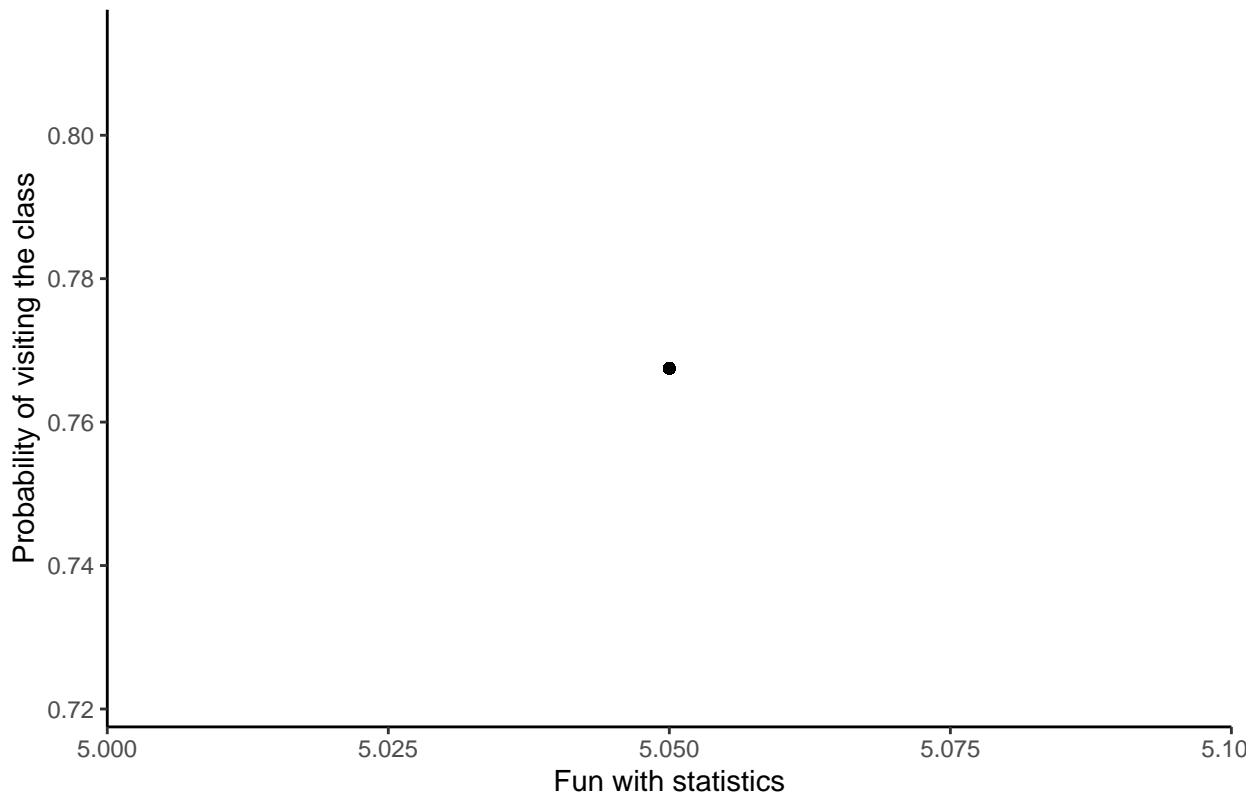
c)

```
# Create a result object
plot1 <- ggplot(stats_lectures, aes(x = stats_joy, y = pred_prob))

# Fill the plot
plot1 <- plot1 + geom_point() +
  labs(x = "Fun with statistics",
       y = "Probability of visiting the class",
       title = "Relation between fun with statistics and visit of statistics classes") +
  theme_classic()

# Show plot
plot1
```

## Relation between fun with statistics and visit of statistics classes



## 4 Model comparison in logistic regression

1. Adapt three more models using the command `glm()`:
  - a) Add predictor `temperature`
  - b) Add predictor `sun_joy`
  - c) Add to the model in b) the interaction of `temperature` and `sun_joy`.
2. Compare the models adapted using the command `anova()`. Which of the models fit the data best, considering parsimony. Use the parameter `test = "Chisq"` to get a p-value.
  - a) Get the `summary()` of the best model.
  - b) Which hypothesis would you deduct from this analysis.
3. Calculate the odds ratios and the confidence intervals for the parameters of the best model.
4. Calculate the pseudo R<sup>2</sup> for the best model (c. f. Field (2012), chapter 8.3.3: Assessing the model: R and R<sup>2</sup>). You may use the function below for this. *Read more on that in Field (2012), R's Souls' Tip 8.2.*
  - a) Copy the code and execute it
  - b) You can apply `logisticPseudoR2s()` to your model.

```
logisticPseudoR2s = function(LogModel) {  
  dev = LogModel$deviance  
  nullDev = LogModel$null.deviance  
  modelN = length(LogModel$fitted.values)  
  R.1 = 1 - dev / nullDev
```

```

R.cs = 1 - exp(-(nullDev - dev) / modelN)
R.n = R.cs / (1 - (exp(-(nullDev / modelN)))))

cat("Pseudo R^2 for logistic regression\n")
cat("Hosmer and Lemeshow R^2 ", round(R.l, 3), "\n")
cat("Cox and Snell R^2 ", round(R.cs, 3), "\n")
cat("Nagelkerke R^2 ", round(R.n, 3), "\n")
}

```

## Solution

```

# the first model
model1 <- glm(attend_or_not ~ stats_joy,
               data = stats_lectures, family = binomial())

# a)
model2 <- glm(attend_or_not ~ stats_joy + temperature,
               data = stats_lectures, family = binomial())

# b)
model3 <- glm(attend_or_not ~ stats_joy + temperature + sun_joy,
               data = stats_lectures, family = binomial())

# c) (we could write that in two manners)
model4 <- glm(attend_or_not ~ stats_joy + temperature*sun_joy,
               data = stats_lectures, family = binomial())

model4 <- glm(attend_or_not ~ stats_joy + temperature +
               sun_joy + temperature:sun_joy,
               data = stats_lectures, family = binomial())

```

## Subtask 1

```
anova(model1, model2, model3, model4, test = "Chisq")
```

## Subtask 2

```

## Analysis of Deviance Table
##
## Model 1: attend_or_not ~ stats_joy
## Model 2: attend_or_not ~ stats_joy + temperature
## Model 3: attend_or_not ~ stats_joy + temperature + sun_joy
## Model 4: attend_or_not ~ stats_joy + temperature + sun_joy + temperature:sun_joy
##   Resid. Df Resid. Dev Df Deviance Pr(>Chi)
## 1       399    433.82
## 2       398    433.20  1     0.622    0.4303
## 3       397    433.12  1     0.083    0.7727

```

```

## 4      396    385.71  1   47.405 5.774e-12 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

It is model 4 that fits the data best.

a)

```

summary(model4)

##
## Call:
## glm(formula = attend_or_not ~ stats_joy + temperature + sun_joy +
##       temperature:sun_joy, family = binomial(), data = stats_lectures)
##
## Deviance Residuals:
##    Min      1Q  Median      3Q     Max
## -2.3668  0.2048  0.6041  0.6903  1.7613
##
## Coefficients: (1 not defined because of singularities)
##                  Estimate Std. Error z value Pr(>|z|)
## (Intercept)      1.43536   0.13866 10.352 < 2e-16 ***
## stats_joy        NA         NA         NA         NA
## temperature      0.05659   0.14162  0.400   0.689
## sun_joy          -0.01291  0.14656 -0.088   0.930
## temperature:sun_joy -0.99906  0.16895 -5.913 3.35e-09 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 433.82 on 399 degrees of freedom
## Residual deviance: 385.71 on 396 degrees of freedom
## AIC: 393.71
##
## Number of Fisher Scoring iterations: 5

```

b) We could deduct the hypothesis:

1. The more students enjoy statistics, the higher is the probability, that they visit statistics classes.
2. Is there an interaction between temperature and enjoyment of sun. Might it be, that higher temperatures lead to less probability to visit classes for students, that enjoy sun?

```

# Odds Ratios
exp(model4$coefficients)

```

### Subtask 3

	(Intercept)	stats_joy	temperature	sun_joy	temperature:sun_joy
##	4.2011631	NA	1.0582247	0.9871727	0.3682260

```
# Confidence intervals, don't forget exp()
exp(confint(model4))
```

```
## Waiting for profiling to be done...

##           2.5 %    97.5 %
## (Intercept) 3.2281763 5.5649819
## stats_joy      NA        NA
## temperature    0.8018372 1.3995432
## sun_joy        0.7398945 1.3166720
## temperature:sun_joy 0.2599905 0.5048007
```

#### Subtask 4

- a) See the text above
- b)

```
logisticPseudoR2s(model4)
```

```
## Pseudo R^2 for logistic regression
## Hosmer and Lemeshow R^2  0.111
## Cox and Snell R^2  0.113
## Nagelkerke R^2  0.171
```

## 5 Render (knit)

Render your document using **Strg + Shift + K** (Windows) oder **Cmd + Shift + K** (Mac). After that, you shoud see a nice looking version of your document. If that works out, congratulations! Your code could be rendered without any errors. If not, don't be frustrated. Fix your errors. We will help you in our exercise class.

## Literature

*Note:* These sheets are based partially on exercises from the book *Discovering Statistics Using R* (Field, Miles & Field, 2012). They've been modified for the porpuses of this seminar, and the R code was updated.

Field, A., Miles, J., & Field, Z. (2012). *Discovering Statistics Using R*. London: SAGE Publications Ltd.

Version: 30 Mai, 2022 21:39