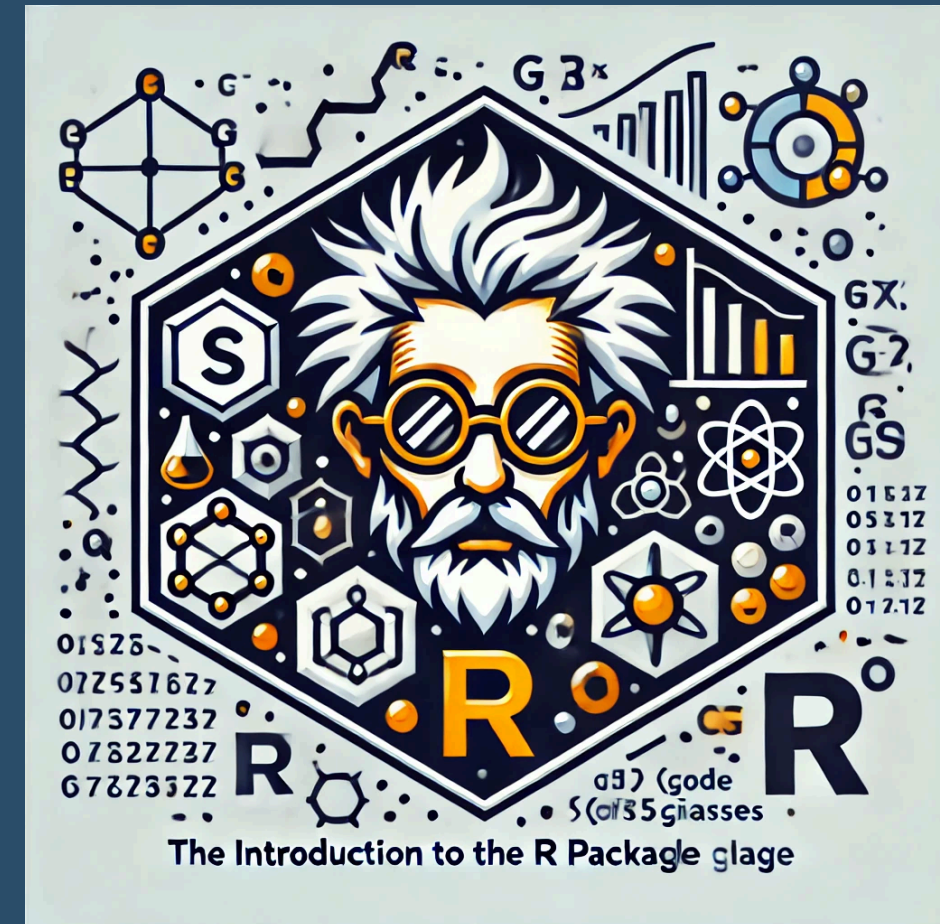


Teil IV – Themes und Data Handling

Grundlagen der Datenanalyse und Statistik mit R | WS 2024/25

Prof. Dr. Daniel D. Schnitzlein



1. Wiederholung aus dem letzten Jahr

Aufgabenblatt III: Datentransformation und Modellierung

(siehe R.Skript auf Moodle und GitHub).

2. Themes: Erstellen von reproduzierbaren publikationsreifen Datenvisualisierungen

Der Ausgangspunkt nach dem zweiten Termin

R Code ↺ Start Over ▶ Run Code

```
1 install.packages('palmerpenguins')
2 library(palmerpenguins)
3 library(ggplot2) # Alternativ könnten wir auch das Meta Paket tidyverse einbinden.
4 library(ggthemes)
```

Warum soll ich ein Theme nutzen?

- Themes haben **keinen** Einfluss darauf, wie die Daten dargestellt oder ob und wie sie transformiert werden.
- Themes ändern die Eigenschaften eines Plots **nicht**, aber sie helfen, den Plot ästhetisch ansprechend zu gestalten oder an einen bestehenden Styleguide anzupassen. Themes steuern Elemente wie z.B. Schriftarten, Ticks, Hintergründe etc.
- Themes vereinfachen es, Abbildungen ein einheitliches Erscheinungsbild innerhalb einer Publikation, innerhalb einer Organisation oder eines Web-Auftritts zu geben.
- Ein einheitliches Erscheinungsbild (das nicht einfach die Standardeinstellungen nutzt) wirkt professioneller und überzeugender verglichen mit wild zusammengewürfelten Designs.
- Individuelle Themes können in **R** vordefiniert werden und müssen daher nur einmal erstellt werden und können dann auf alle Arten von Abbildungen angewendet werden.

Kernelemente des **ggplot2**-Theme-Systems

elements	Definiert das Objekt, das verändert werden soll.	z.B. <code>plot.title</code> oder <code>axis.ticks.x</code> usw. Eine vollständige Liste aller elements findet sich hier .
element function	Jedes element ist mit einer oder mehreren element functions verbunden.	z.B. <code>element_text(...)</code> definiert die Schriftgröße, Schriftart, Farbe etc. von <code>plot.title</code> .
theme(...)	Funktion um die Änderungen durchzuführen.	z.B. <code>theme(plot.title = element_text(colour = "red"))</code>
vollständige Themes	Vollständige Themes , die entweder bereits vordefiniert oder selbst definiert werden können.	z.B. <code>theme_bw()</code>

Vollständige **ggplot2**-Themes (1)

ggplot2 ermöglicht es anhand von Themes die Formatierung von Abbildungen weitgehend zu automatisieren (z.B. Anwendung von einheitlichem Corporate Design).

R Code

↺ Start Over

▶ Run Code

```
1 ggplot(data = penguins, aes(x = flipper_length_mm, y = body_mass_g)) +  
2   geom_point(aes(color = species, shape = species)) +  
3   geom_smooth(method = "lm") +  
4   labs(title = "Body mass and flipper length",  
5         subtitle = "Dimensions for Adelie, Chinstrap, and Gentoo Penguins",  
6         x = "Flipper length (mm)",  
7         y = "Body mass (g)",  
8         color = "Species",  
9         shape = "Species") +  
10  theme_economist()
```


Vollständige ggplot2-Themes (2)

<code>theme_gray()</code>	Das charakteristische ggplot2-Design mit grauem Hintergrund und weißen Gitterlinien, das die Daten hervorhebt und Vergleiche erleichtert.
<code>theme_bw()</code>	Variation des klassischen „dark-on-light“-ggplot2-Themes. Eignet sich möglicherweise besser für die Anzeige mit einem Beamer.
<code>theme_linedraw()</code>	Theme nutzt nur schwarzen Linien unterschiedlicher Breite auf weißem Hintergrund.
<code>theme_light()</code>	Ein ähnliches Theme wie <code>theme_linedraw()</code> , aber mit hellgrauen Linien und Achsen, um mehr Aufmerksamkeit auf die Daten zu lenken.
<code>theme_dark()</code>	Der dunkle Cousin von <code>theme_light()</code> , mit ähnlichen Strichstärken, aber dunklem Hintergrund.
<code>theme_minimal()</code>	Ein minimalistisches Theme ohne Hintergrund.
<code>theme_classic()</code>	Ein klassisch aussehendes Theme mit x- und y-Achsenlinien und ohne Gitterlinien.
<code>theme_void()</code>	Ein komplett leeres Theme.
<code>theme_test()</code>	Ein Theme zum Testen visueller Komponenten. Es sollte sich idealerweise

Probieren Sie es aus!

R Code

↺ Start Over

▶ Run Code

```
1 ggplot(data = penguins, aes(x = flipper_length_mm, y = body_mass_g)) +  
2   geom_point(aes(color = species, shape = species)) +  
3   geom_smooth(method = "lm") +  
4   labs(title = "Body mass and flipper length",  
5         subtitle = "Dimensions for Adelie, Chinstrap, and Gentoo Penguins",  
6         x = "Flipper length (mm)",  
7         y = "Body mass (g)",  
8         color = "Species",  
9         shape = "Species") +  
10  theme_gray()
```

Anpassen der integrierten themes

- Alle vordefinierten Themes haben einen Parameter `base_size()` mit dem die Basisschriftgröße (Achsentitel) aller Schriften mit einem Befehl geändert werden kann.
- Dabei wird auch die Achsenbeschriftung (0.8x) und der Titel (1.2x) angepasst.
- Mit `set_theme()` kann ein bestimmtes theme als Standardtheme gesetzt werden.
- Alle weitergehenden Änderungen müssen über die Ansprache der entsprechenden Elemente direkt erfolgen.

Beispiel: Änderung der Farbe im Titel

R Code

↻ Start Over

▶ Run Code

```
1 ggplot(data = penguins, aes(x = flipper_length_mm, y = body_mass_g)) +  
2   geom_point(aes(color = species, shape = species)) +  
3   geom_smooth(method = "lm") +  
4   labs(title = "Body mass and flipper length",  
5         subtitle = "Dimensions for Adelie, Chinstrap, and Gentoo Penguins",  
6         x = "Flipper length (mm)",  
7         y = "Body mass (g)",  
8         color = "Species",  
9         shape = "Species") +  
10  theme_bw() +  
11  theme(plot.title = element_text(colour = "red"))
```

Beispiel: Volle Individualisierung

R Code

↺ Start Over

▶ Run Code

```
1 ggplot(data = penguins, aes(x = flipper_length_mm, y = body_mass_g)) +
2   geom_point(aes(color = species, shape = species)) +
3   geom_smooth(method = "lm") +
4   labs(title = "Body mass and flipper length",
5         subtitle = "Dimensions for Adelie, Chinstrap, and Gentoo Penguins",
6         x = "Flipper length (mm)",
7         y = "Body mass (g)",
8         color = "Species",
9         shape = "Species") +
10  theme_bw() +
11    theme(plot.title = element_text(colour = "red", family="Comic Sans MS"),
12          plot.subtitle = element_text(colour = "blue", face="bold"),
13          plot.title.position = "plot",
14          axis.title = element_text(colour = "blue"),
15          axis.text = element_text(colour = "purple"),
16          legend.position = "top",
17          legend.location = "plot",
18          legend.margin = margin(0, 0, 0, 0), # turned off for alignment
19          legend.justification.top = "left",
20          legend.justification.left = "top",
21          panel.grid.major = element_blank(),
22          panel.grid.minor = element_blank())
```

on GitHub Repo: https://github.com/dschnitzlein/GDSR_4_Part_1/

```
22     panel.grid.minor = element_blank(),  
23     panel.border=element_blank(),  
24     axis.line.x = element_line(colour = "lightgray"),  
25     axis.line.y = element_line(colour = "lightgray")  
26 }
```

Definieren eigener Themes

R Code

↺ Start Over

▶ Run Code

```
1 Daniels_theme <- function() {  
2   theme_bw() +  
3     theme(plot.title = element_text(colour = "red", family="Comic Sans MS"),  
4           plot.subtitle = element_text(colour = "blue", face="bold"),  
5           plot.title.position = "plot",  
6           axis.title = element_text(colour = "blue"),  
7           axis.text = element_text(colour = "purple"),  
8           legend.position = "top",  
9           legend.location = "plot",  
10          legend.margin = margin(0, 0, 0, 0), # turned off for alignment  
11          legend.justification.top = "left",  
12          legend.justification.left = "top",  
13          panel.grid.major = element_blank(),  
14          panel.grid.minor = element_blank(),  
15          panel.border=element_blank(),  
16          axis.line.x = element_line(colour = "lightgray"),  
17          axis.line.y = element_line(colour = "lightgray")  
18    )  
19 }
```

Nutzen eigener Themes

R Code

↺ Start Over

▶ Run Code

```
1  
2 theme_set(Daniels_theme())  
3  
4 ggplot(data = penguins, aes(x = flipper_length_mm, y = body_mass_g)) +  
5   geom_point(aes(color = species, shape = species)) +  
6   geom_smooth(method = "lm") +  
7   labs(title = "Body mass and flipper length",  
8         subtitle = "Dimensions for Adelie, Chinstrap, and Gentoo Penguins",  
9         x = "Flipper length (mm)",  
10        y = "Body mass (g)",  
11        color = "Species",  
12        shape = "Species")
```


Wichtige Themeelemente: 1. Plotelemente

Element	Funktion	Beschreibung
plot.background	<code>element_rect()</code>	Hintergrund der Abbildung (ohne Daten)
plot.title	<code>element_text()</code>	Titel der Abbildung
plot.margin	<code>margin()</code>	Margins um die Abbildung
plot.title.position	"panel", "plot"	Alignment des Titels entweder mit dem Panel oder Plot

Wichtige Themeelemente: 2. Achsenelemente

Element	Funktion	Beschreibung
axis.line	<code>element_line()</code>	Achsenlinie (in std. Themes versteckt)
axis.text	<code>element_text()</code>	Labels der Achsenticks
axis.text.x	<code>element_text()</code>	x-Achsentickslabels
axis.text.y	<code>element_text()</code>	y-Achsentickslabels
axis.title	<code>element_text()</code>	Achsentitel
axis.title.x	<code>element_text()</code>	Achsentitel x-Achse
axis.title.y	<code>element_text()</code>	Achsentitel y-Achse
axis.ticks	<code>element_line()</code>	Achsentickstriche
axis.ticks.length	<code>unit()</code>	Achsentickstrichlänge

Wichtige Themeelemente: 3. Legendenelemente

Element	Funktion	Beschreibung
legend.background	<code>element_rect()</code>	Legendenhintergrund
legend.key	<code>element_rect()</code>	Hintergrund der Legendeneinträge
legend.key.size	<code>unit()</code>	Größe der Legendeneinträge
legend.key.height	<code>unit()</code>	Höhe der Legendeneinträge
legend.key.width	<code>unit()</code>	Breite der Legendeneinträge
legend.margin	<code>unit()</code>	Margin um die Legende
legend.text	<code>element_text()</code>	Legenden Labels
legend.text.align	0–1	Legenden Label Alignment (0 = rechts, 1 = links)
legend.title	<code>element_text()</code>	Legendenname
legend.title.align	0–1	Legendname Alignment (0 = rechts, 1 = links)
legend.position	"top", "bottom", "left", "inside", "right", "none"	Positionieren der Legende
legend.justification.top	"left", "center", "right"	Positionieren der Legende (ebenfalls für alle Pos.)

Wichtige Themeelemente: 4. Panelelemente

Element	Funktion	Beschreibung
panel.background	<code>element_rect()</code>	Panel Hintergrund (unter den Daten)
panel.border	<code>element_rect()</code>	Panelrahmen
panel.grid.major	<code>element_line()</code>	Hauptgridlines
panel.grid.major.x	<code>element_line()</code>	Hauptgridlines vertikal
panel.grid.major.y	<code>element_line()</code>	Hauptgridlines horizontal
panel.grid.minor	<code>element_line()</code>	Minorgridlines
panel.grid.minor.x	<code>element_line()</code>	Minorgridlines vertikal
panel.grid.minor.y	<code>element_line()</code>	Minorgridlines horizontal
aspect.ratio	numeric	Plot aspect ratio

Wichtige Themeelemente: 5. Facetingelemente

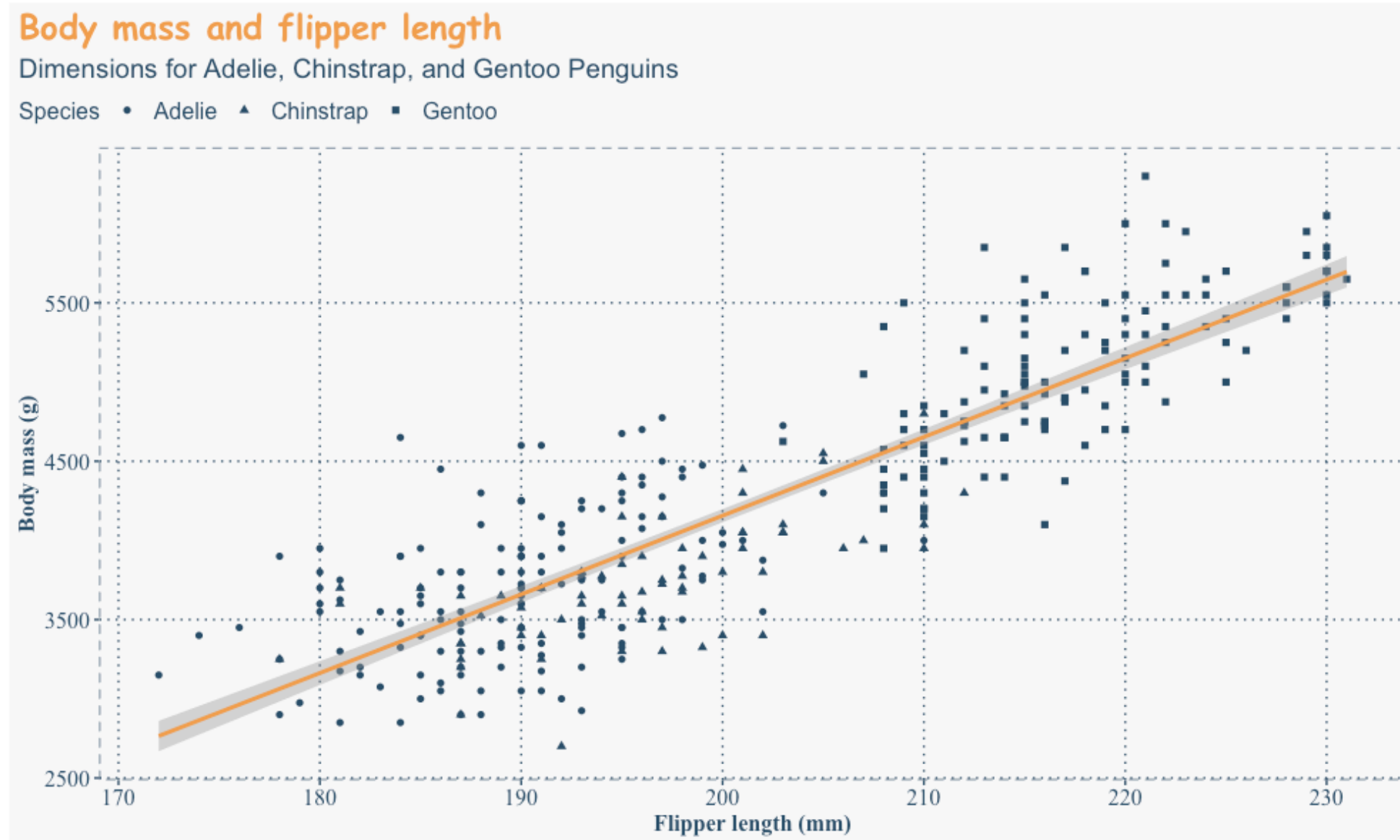
Element	Funktion	Beschreibung
strip.background	<code>element_rect()</code>	Hintergrund der Panel Strips
strip.text	<code>element_text()</code>	Strip Text
strip.text.x	<code>element_text()</code>	Horizontaler Strip Text
strip.text.y	<code>element_text()</code>	Vertikaler Strip Text
panel.spacing	<code>unit()</code>	Margin zwischen Facets
panel.spacing.x	<code>unit()</code>	Margin zwischen Facets (vertikal)
panel.spacing.y	<code>unit()</code>	Margin zwischen Facets (horizontal)

Wichtige Theme-Funktionen

<code>element_text(...)</code>	Kontrolliert Schriftart (<code>family</code>), Typ (<code>face</code>), Farbe (<code>colour</code>), Größe (<code>size</code>) usw.
<code>margin(0,0,0,0)</code>	Kontrolliert den Rand um ein Objekt
<code>element_line(...)</code>	Kontrolliert die Farbe (<code>colour</code>), Linienart (<code>linetype</code>) und Dicke (<code>linewidth</code>) einer Linie
<code>element_rect(...)</code>	Kontrolliert Rechtecke mit den Parametern <code>fill</code> , <code>linewidth</code> , <code>colour</code> und <code>linetype</code>
<code>element_blank(...)</code>	Zeichnet nichts (zum Deaktivieren von Elementen)

Aufgabe 1

Erstellen Sie ihr eigenes Theme um die Abbildung folgendermaßen aussehen zu lassen:



GitHub-Repo: https://github.com/dschnitzlein/GDSR_4_Part_1/

Aufgabe 1 – Teil I

Erstellen Sie bitte ein eigenes Theme mit dem Namen `GDSR_theme`, das auf dem `theme_bw()` basiert und folgendem Design folgt:

- Die gesamte Abbildung (Plot, Legende, und Panel) soll einen leicht angegrauten Hintergrund haben: `#fafafa`
- Alle Schriftelemente mit Ausnahme des Titels sollen ein dunkles Blau als Farbe haben: `#2A4D69`
- Alle Schriftelemente mit Ausnahme des Titels (Achsenbeschriftung, Achsentitel, Achsenticks, Legendentext, Legendentitel, Subtitle) sollen in **Times New Roman** gesetzt werden.
- Der Titel soll in **Comic Sans MS** gesetzt werden und folgende Farbe aus dem Design haben: `#F29E4C`
- Achsentitel und Abbildungstitel sollen fett gedruckt sein.
- Legendentext und Legendentitel sowie Achsentext und Achsentitel haben Schriftgröße 12. Der Titel hat Schriftgröße 18 und der Untertitel hat Schriftgröße 14.
- Das Panel enthält nur Major-Gridlines in Farbe `#2A4D69` und gestrichelt (`linetype 3`) in `size = 0.5`.
- Die Abbildung hat einen Rahmen (`panel.border`) ebenfalls in Farbe `#2A4D69` und mit `linetype 2`.
- Die Achsenticks sind ebenfalls in Farbe `#2A4D69`
- Die Legende soll oben links platziert sein.
- Titel und Legende sollen linksbündig mit der Achsenbeschriftung ausgerichtet werden.

Aufgabe 1 – Teil II

R Code ↺ Start Over ▶ Run Code

```
1 library(palmerpenguins)
2 library(ggplot2)
3 ggplot(data = penguins, aes(x = flipper_length_mm, y = body_mass_g)) +
4   geom_point(aes(color = species, shape = species)) +
5   geom_smooth(method = "lm") +
6   labs(title = "Body mass and flipper length",
7         subtitle = "Dimensions for Adelie, Chinstrap, and Gentoo Penguins",
8         x = "Flipper length (mm)",
9         y = "Body mass (g)",
10        color = "Species",
11        shape = "Species")
```

Ändern Sie nun folgende Punkte an der Abbildung:

- Unterscheiden Sie die Arten von Pinguinen nur durch Shape und nicht durch Color.
- Färben Sie die Datenpunkte einheitlich in #2A4D69 ein.
- Ändern Sie die Farbe der Regressionslinie zu #F29E4C.
- Weisen Sie der Abbildung das neue Theme zu.

Speichern der Abbildungen

```
R Code ↺ Start Over ▶ Run Code
1 ggsave(
2   filename,
3   plot = last_plot(),
4   device = NULL,
5   path = NULL,
6   scale = 1,
7   width = NA,
8   height = NA,
9   units = c("in", "cm", "mm", "px"),
10  dpi = 300,
11  limitsize = TRUE,
12  bg = NULL,
13  create.dir = FALSE,
14  ...
15 )
```

`ggsave()` speichert in viele Formate, z.B. "png", "eps", "ps", "tex" (pictex), "pdf", "jpeg", "tiff", "png", "bmp", "svg", ... Neben der Größe (Breite und Höhe) kann auch die Auflösung (dpi) oder das Farbmodel angegeben werden. Anmerkung: Die `base-R` Lösung über `pdf()` oder `png()` ist auch mit `ggplot`-Objekten möglich.

3. Einlesen und Bereinigen von Daten

Das Menü

- Die meisten Datensätze können interaktiv über das Menü in R-Studio eingelesen werden.
- Dabei wird auch der entsprechende R-Code generiert, den Sie dann in Ihr Skript übernehmen können.
- Aufgrund der einfachen Bedienung und der komfortablen Vorschau-Funktion ist es empfehlenswert diesen Weg für das erste Mal des Einlesens eines Datensatzes zu wählen.
- Natürlich können alle Datensätze auch über ein Skript eingelesen werden.
- Bei spezielleren Datensätzen ist dieser Weg zu bevorzugen.
- Wir beschränken uns hier im Kurs auf "Rechtecksdaten" in den Formaten *.csv und *.xls bzw. *.xlsx

CSV-Datensätze (1)

R Code ↺ Start Over ▶ Run Code

```
1 library(tidyverse)
2 students <- read_csv("data/students.csv")
```

alternativ

R Code ↺ Start Over ▶ Run Code

```
1 students_web <- read_csv("https://pos.it/r4ds-students-csv")
```

Clean-up missing values

R Code ↺ Start Over ▶ Run Code

```
1 students <- read_csv("data/students.csv")
```

alternativ

R Code ↺ Start Over ▶ Run Code

```
1 students <- read_csv("data/students.csv", na=c("N/A", ""))
```

Einheitliche Namen der Variablen

R Code ↺ Start Over ▶ Run Code

```
1 install.packages("janitor")
2 library(janitor)
3
4 students <- students |>
5   clean_names()
```

Variablentypen zuweisen

Die Variable `meal_plan` ist eine kategoriale Variable die in R als Faktor angelegt sein sollte. Nur so kann R später für Modellierung etc. die Variablen korrekt verwenden.

R Code ↺ Start Over ▶ Run Code

```
1
2 students <- students |>
3   mutate(meal_plan = factor(meal_plan))
4
5 students
```


Variablen bereinigen

Die Variable **age** enthält aktuell Zahlen und Strings. Das wollen wir bereinigen.

R Code ↺ Start Over ▶ Run Code

```
1
2 students <- students |>
3   mutate(
4     age = parse_number(if_else(age == "five", "5", age))
5   )
6
7 students
```

Weitere Informationen und Optionen zum Umgang mit CSV-Files finden Sie unter [?read_csv](#).

Excel-Daten

- Microsoft Excel ist wahrscheinlich das meistgenutzte Tabellenkalkulationsprogramm und wahrscheinlich auch das meistgenutzte Programm zur (oberflächlichen) Datenanalyse.
- Die Excel-Datenformate sind daher auch häufig genutzte Output-Formate für Datengenerierungen und Datenweitergaben (speziell in Unternehmen).
- R kann Excel-Daten u.a. über das Paket `readxl` lesen und via `writexl` schreiben. Die zugehörige Funktion ist `read_excel()`.
- Anmerkung: R kann auch mit Stata, SPSS, SAS o.ä. Daten umgehen z.B. via des Pakets `haven`.

Excel-Daten einlesen (1)

R Code ↺ Start Over ▶ Run Code

```
1 library(readxl)
2
3 students <- read_excel("data/students.xlsx")
```

Excel-Daten einlesen (2)

R Code ↺ Start Over ▶ Run Code

```
1  
2 students <- read_excel("data/students.xlsx",  
3   col_names = c("student_id", "full_name", "favourite_food", "meal_plan", "age"),  
4   skip = 1  
5 )
```

Excel-Daten aus Dateien mit mehreren Tabellenblättern einlesen

R Code ↺ Start Over ▶ Run Code

```
1
2 penguins_torgersen <- read_excel("data/penguins.xlsx", sheet = "Torgersen Island", na = "NA")
3 penguins_biscoe <- read_excel("data/penguins.xlsx", sheet = "Biscoe Island", na = "NA")
4 penguins_dream <- read_excel("data/penguins.xlsx", sheet = "Dream Island", na = "NA")
```

Daten zu einem Datensatz zusammenfügen

R Code ↺ Start Over ▶ Run Code

```
1 penguins <- bind_rows(penguins_torgersen, penguins_biscoe, penguins_dream)
2 penguins
```

Unsere Themenliste

1. Einführung in R und R-Studio: Überblick über die Installations- und Einrichtungsprozesse | Grundlegende Funktionen und Bedienung von R und RStudio
2. Grundlagen der Statistiksprache R: Syntax und Datenstrukturen in R | Einführung in Funktionen und Pakete
3. Datenmanagement in R: Methoden der Datenorganisation und -vorbereitung | Importieren, Bereinigen und Transformieren von Datensätzen
4. Einführung in die Pakete des tidyverse: Überblick über die wichtigsten tidyverse-Pakete wie z.B. dplyr und ggplot2 | Anwendung dieser Pakete zur effizienten Datenanalyse und -visualisierung
5. Deskriptive Statistik in R: Berechnung und Interpretation grundlegender statistischer Kennzahlen | Anwendung von deskriptiven Methoden zur Datenexploration | Einführung in die statistische Modellierung am Beispiel linearer Modelle
6. Datenvisualisierung in R: Erstellen von publikationsreifen Grafiken und Diagrammen mit ggplot2 | Gestaltung und Interpretation von Datenvisualisierungen zur Unterstützung der Datenanalyse

Noch offen: Einheitliche und publikationsreife Tabellen (nächste Woche).