

# Statistik und Datenanalyse: Aufbau

*Mittelwert, Standardabweichung, Korrelation*

Benjamin Fretwurst

▶ PDF-Version der Folien



# Inhalt

- 1 Modulintro
  - 1.1 Lernziele des Moduls
  - 1.2 Leistungsnachweis
  - 1.3 Inhalt und Aufwand
  - 1.4 Material
- 2 Befragung
- 3 Was bisher geschah
  - 3.1 Univariate Statistik
  - 3.2 Bivariate Statistik
  - 3.3 Hypothesentesten
- LEF 1
  - Essayfragen
  - MC-Fragen
- Take Home
- Ausblick



# 1 Modulintro



# Lernziele heute

## Modulorga

- Lernziele und Inhalte von «Statistik: Aufbau»
- Orga
- Lernmaterialien

## Anschluss an Statistik: Einführung

- Mittelwert
- Streuung
- Inferenz
  - Punktschätzung (Standardfehler)
  - Intervallschätzung
  - Tests
- Kovarianz und Korrelation



# 1.1 Lernziele des Moduls

## Grundlagenkenntnisse

- Sie erlangen Kenntnisse multivariater Statistik.
- Verständnis empirischer Forschungsbeiträge
- Umgang mit Vorlagen für die Anwendung der Verfahren in R.

## Praktische Statistik für «Methoden Aufbau»

Sie können Analysemethoden anwenden, die Sie für Ihr Projekt in Methoden Aufbau brauchen.

## Verständnis fortgeschrittener Statistik für die Forschungsseminare

Vorbereitung auf die Forschungsseminare, in denen Sie weiterführende Analyseverfahren verstehen müssen (Texte und Analysen) und teilweise auch anwenden.



# 1.1.1 Voraussetzungen

## Inhaltliche Voraussetzungen

- Deskriptive und induktive Statistik aus «Statistik und Datenanalyse: Einführung»
- heute Rückblick
- Folien der Statistik-Einführungs-Vorlesung (+ Buch Ihrer Wahl)

## Technische Voraussetzung

- R
- R-Studio (oder VS-Code)
- Weiteres siehe Anleitung [R-Installationsanleitung](#)



# 1.2 Leistungsnachweis

## Was ist Prüfungstoff?

- Folien
- Vorlesung
- Begleittext

## Prüfungstermin (Major + Minor)

- Hauptklausur
  - 03.01.2024
  - 14:30–15:30 (60 Minuten)
  - BYOD ⇨ [Y-15-G-60 + Y24-G-45]
- Wiederholungsklausur Major und Minor
  - 07.02.2024
  - 14:30–15:30 (60 Minuten)
  - BYOD ⇨ Y-15-G-60



## 1.3 Inhalt und Aufwand





# Inhalte

1. Uni- und Bivariate Statistik
2. GLM – Regression
3. GLM – BLUE
4. Übung: GLM I
5. GLM – Kategoriale UV
6. GLM – Interaktionen
7. GLM – Übung II
8. Dimensionsreduktion
9. Übung: Dimensionsreduktion
10. LogReg und ML
11. Übung: Machine Learning
12. Clusteranalyse
13. Übung: R
14. Zusammenfassung



## Aufwandt für 6 ECTS

Aufwand	in h	h/Woche	Punkte
Besuch der Vorlesung	21	1.5	0.7
Vor und Nachbereitung	21	1.5	0.7
Lesen der Texte	42	3.0	1.4
Übungsaufgaben in R	50	3.6	1.7
Prüfungsvorbereitung	42	3.0	1.4
Studienteilnahmepunkte	4	0.3	0.1
Summe für 6 ECTS * 30h	180	12.9	6.0




# 1.4 Material



# E-Learning – OLAT

## OLAT

- Mitteilungen
- Forum
  - Diskussionen und Fragen
  - Klausurvorbereitung
- [Studienteilnahmepunkte](#) 



# R-Seite und Begleittext (BETA!)

## R-Seite

Auf [r.ikmz.uzh.ch/Wissen\\_macht\\_R!](http://r.ikmz.uzh.ch/Wissen_macht_R!) finden Sie Anleitungen zur Installation von R und R-Studio sowie Beispielskripte und Vorlagen, die wir vor allem auch für brauchen, bzw. Ihnen helfen sollen.

## Begleittext ist noch BETA!

Der Begleittext auf [stat.ikmz.uzh.ch/Aufbau](http://stat.ikmz.uzh.ch/Aufbau) wird während des Semesters deutlich überarbeitet und ist daher BETA! Es gibt keine vollständige Gewähr für Fehlerfreiheit!



# Zusatzliteratur (Wiederholung Statistik Einführung)

Für die Wiederholung von  
«Statistik: Einführung».

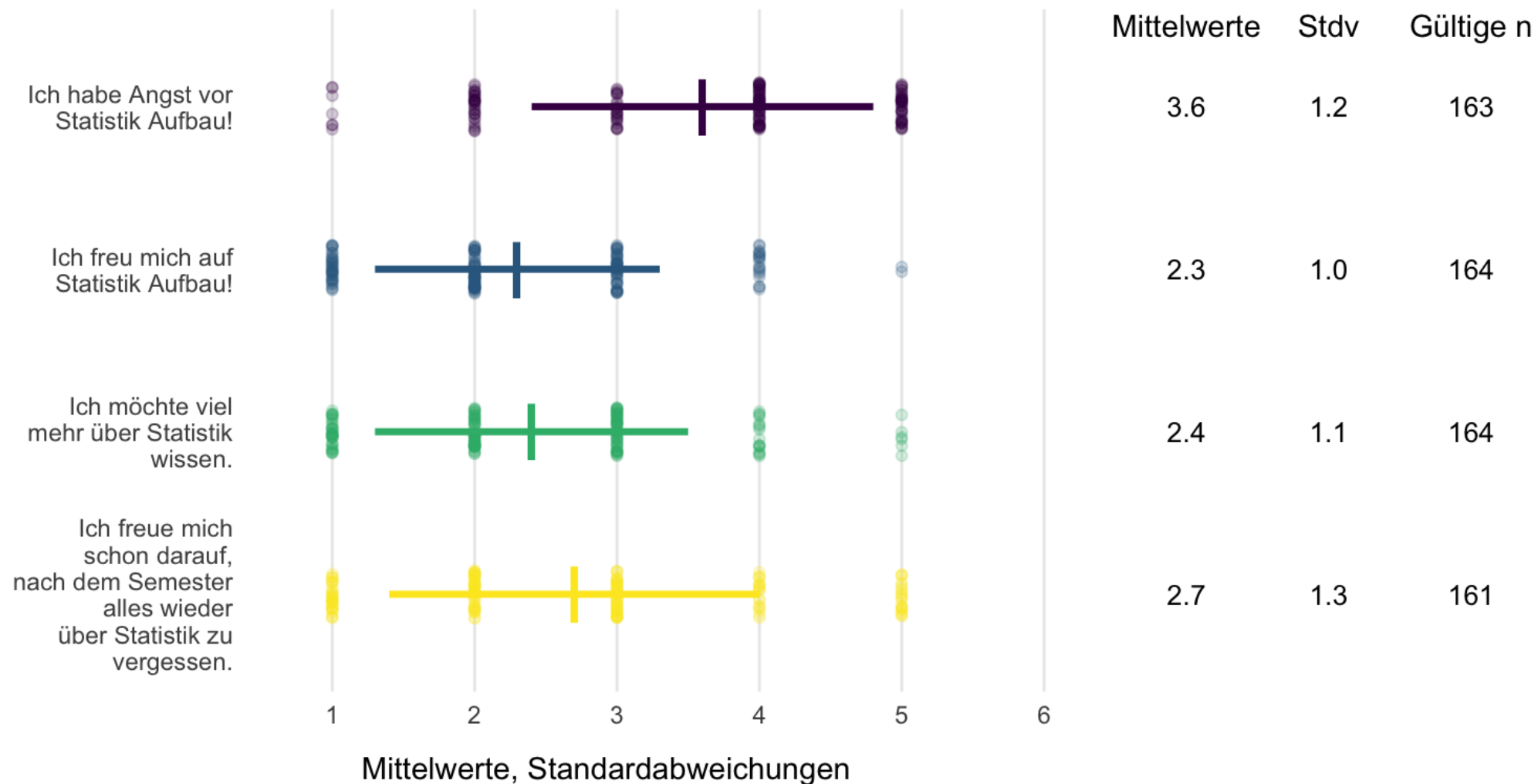


# 2 Befragung



# Erwartungen

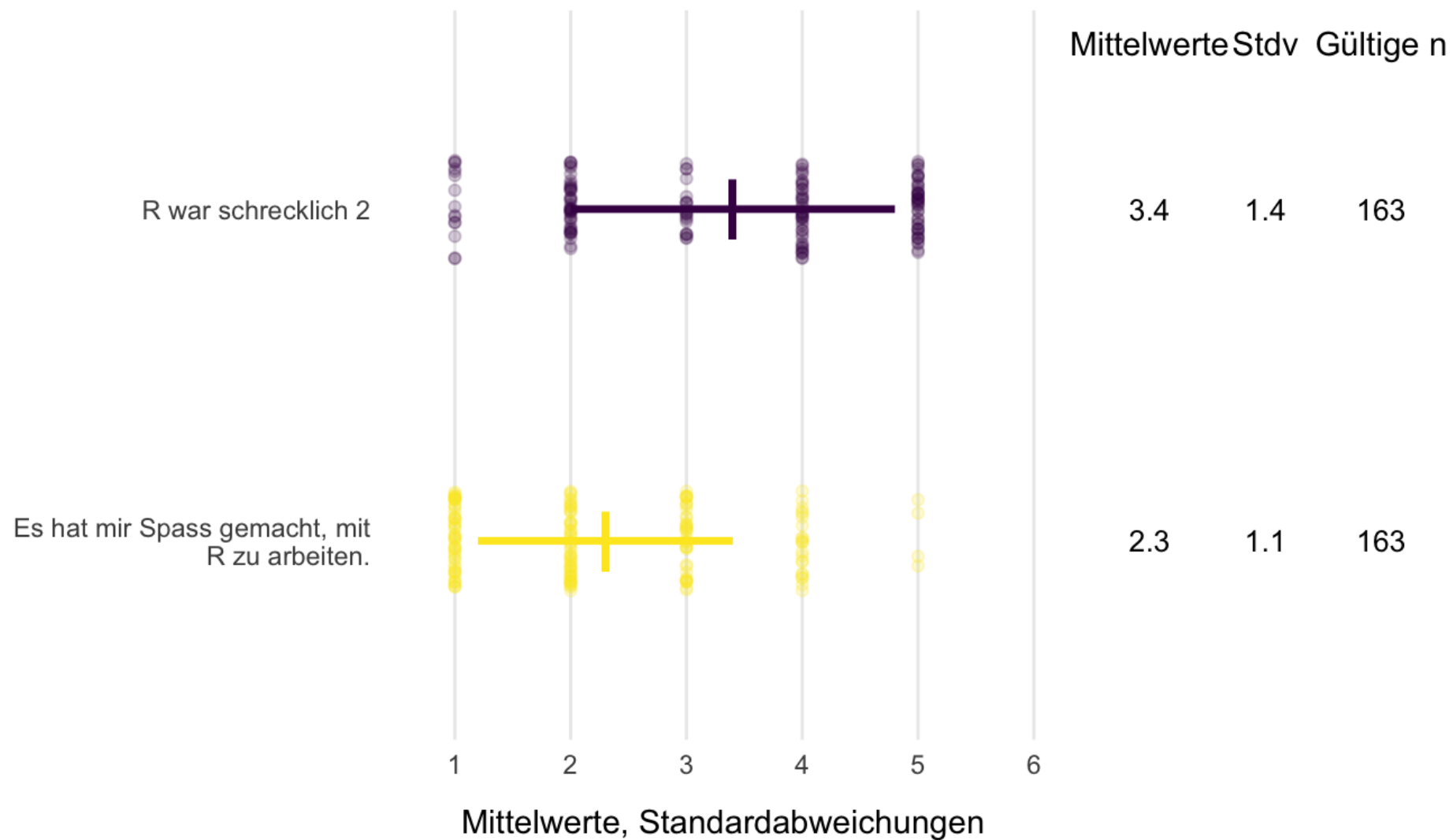
## ► R-Code anzeigen





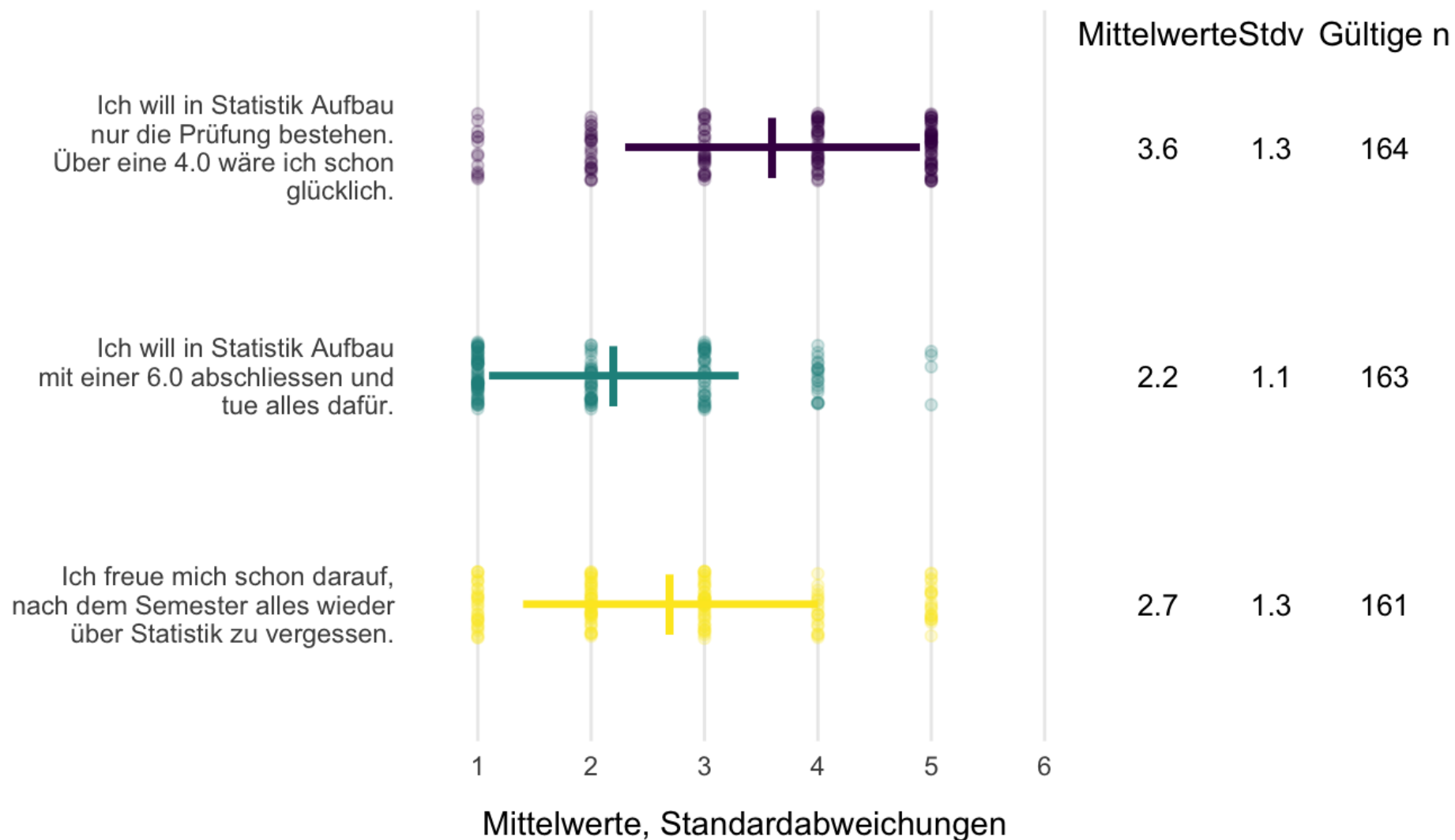
# Umgang mit R

## ► R-Code anzeigen



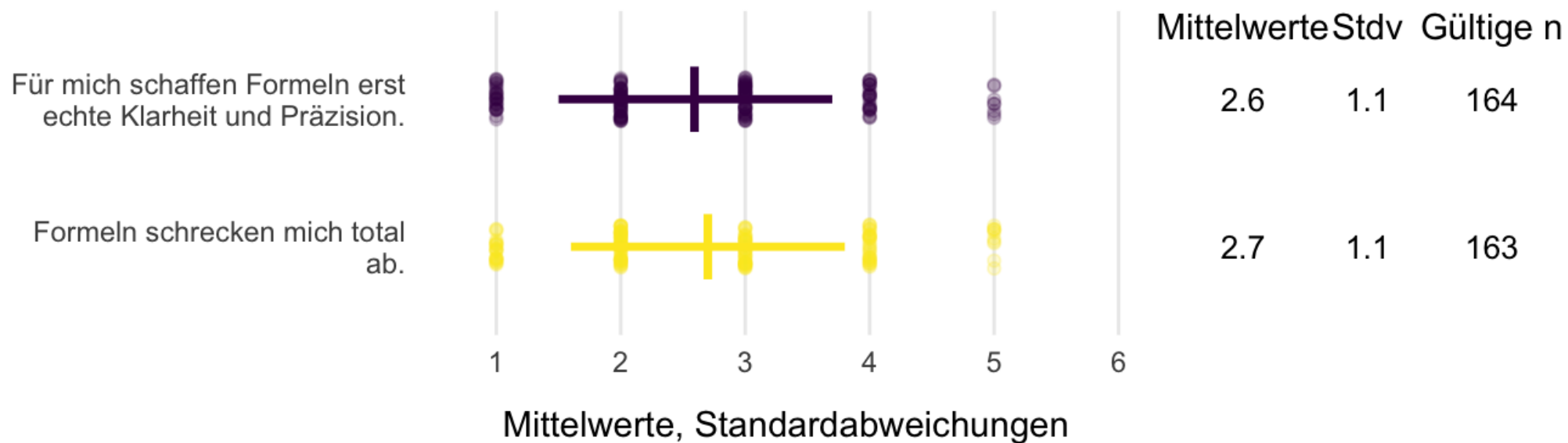
# Ziele

## ► R-Code anzeigen



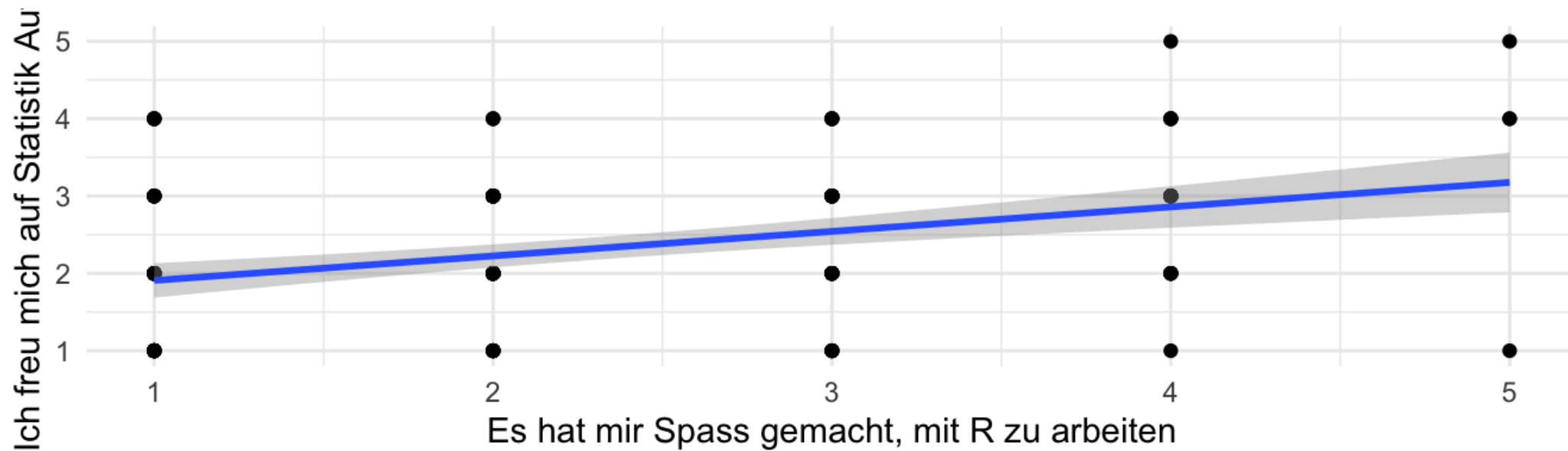
# Was geht mit Formeln?

## ► R-Code anzeigen



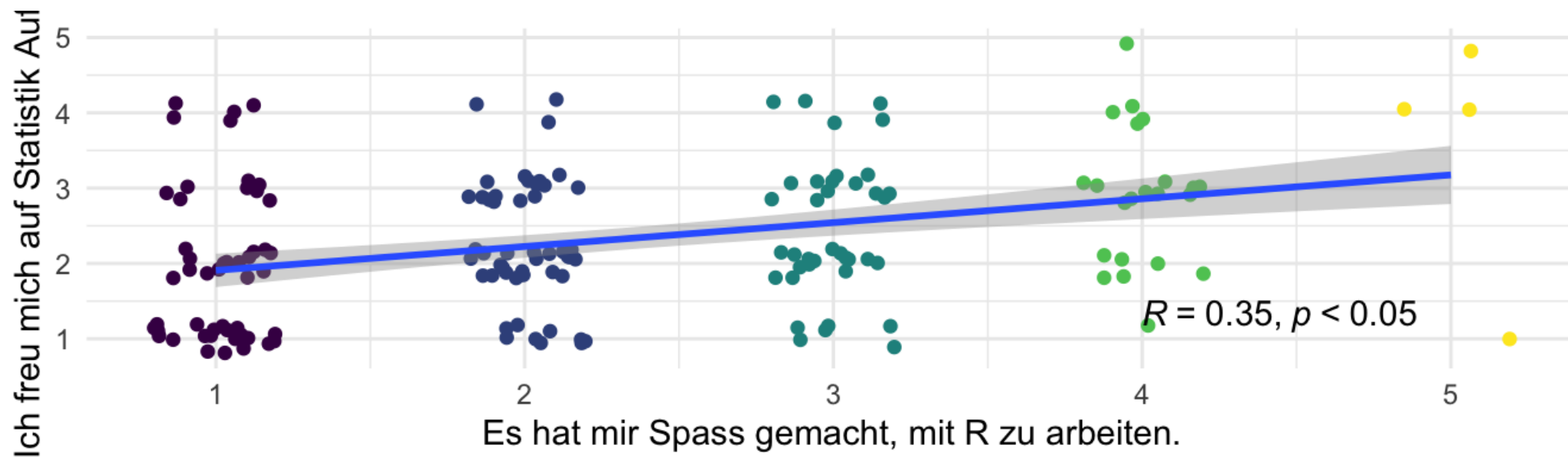
# Spas und Freude

## ► R-Code anzeigen



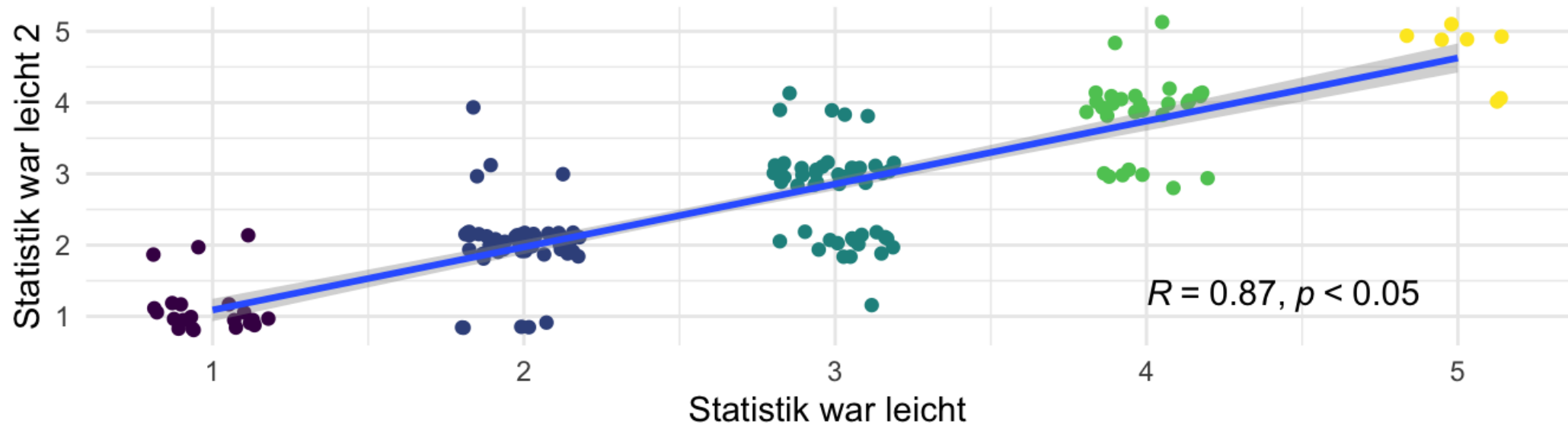
# Statistik war leicht

## ► R-Code anzeigen



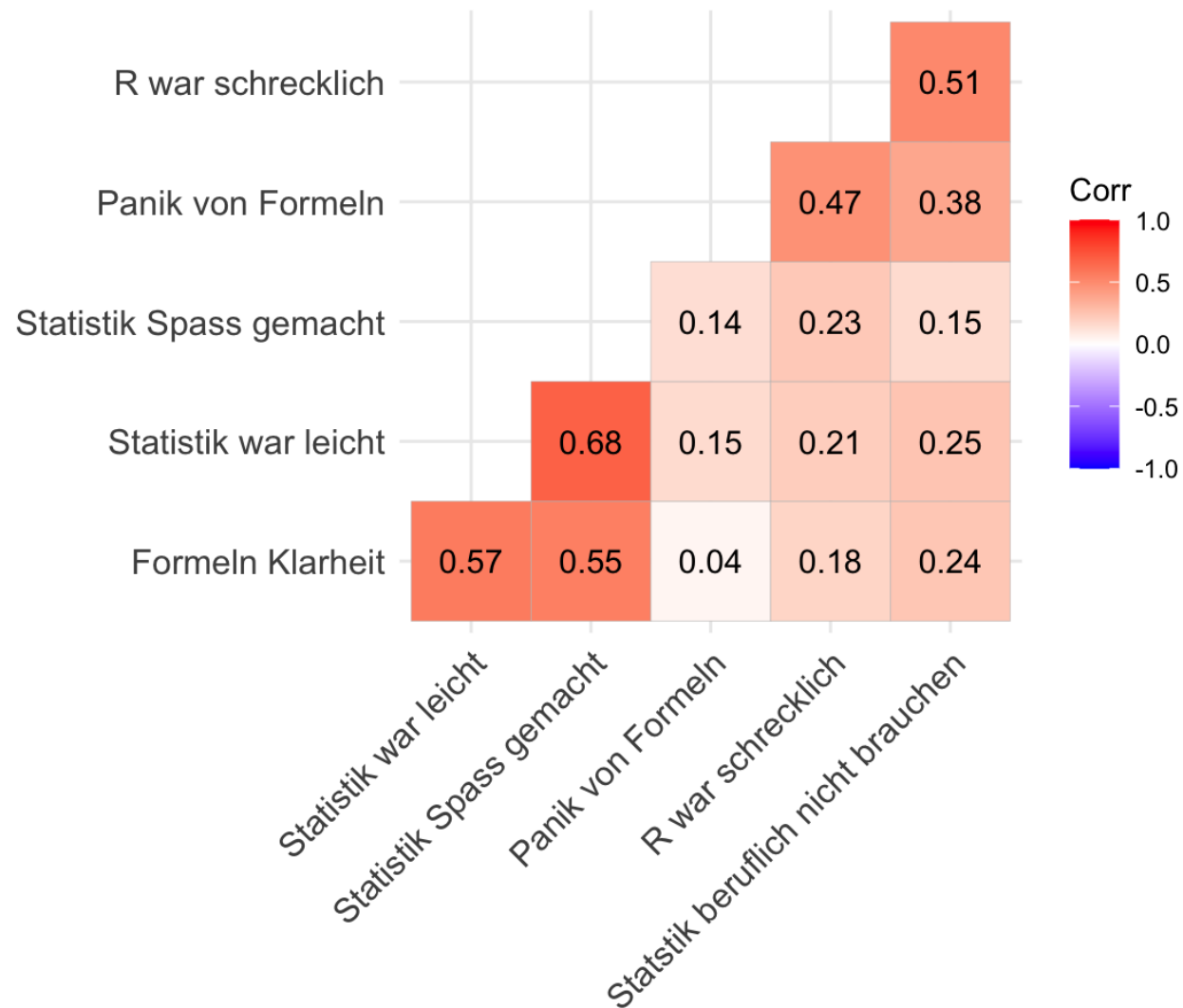
# Statistik war leicht (doppelt)

## ► R-Code anzeigen



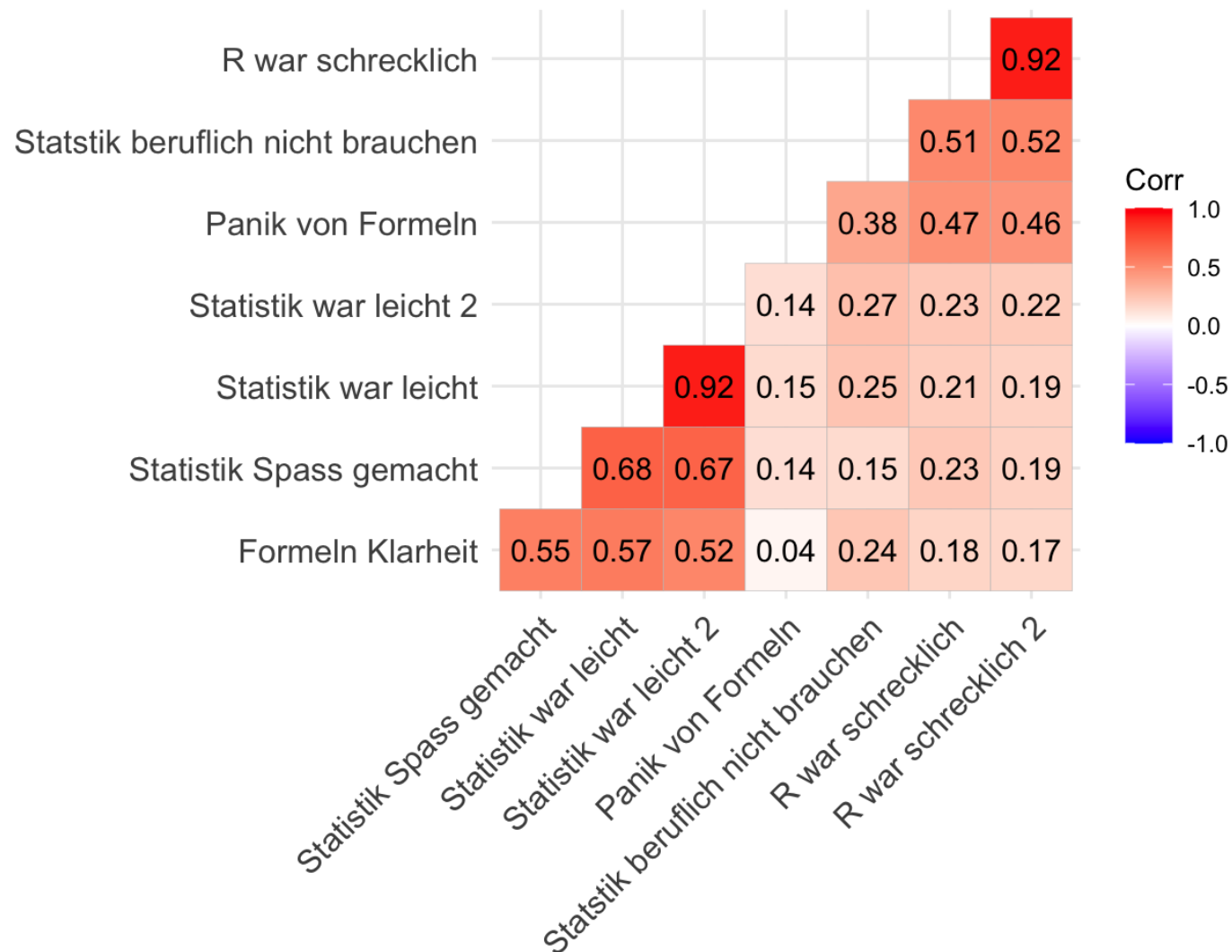
# Korrelationen

## ► R-Code anzeigen



# Korrelation doppelte Frage

## ► R-Code anzeigen





# 3 Was bisher geschah



# 3.1 Univariate Statistik



## 3.1.1 Mittelwert

Der Mittelwert ist das Gleiche wie ein Durchschnitt.

$$\bar{x} = \frac{1}{n} \sum_i^n (x_i)$$

$$\bar{y} = \frac{1}{n} \sum_i^n (y_i)$$

Mittelwert wird geschrieben als Kennwert  $\bar{x}$  und Parameter  $\mu$ .



## 3.1.2 Varianz

Die Varianz ist der Mittelwert (Durchschnitt) der quadrierten Abweichungen vom Mittelwert.

$$\sigma^2 = V = \frac{1}{n} \sum_i^n (x_i - \mu)^2$$

Die Varianz als unbekannter Parameter wird als  $\sigma^2$  gekennzeichnet.



### 3.1.3 Standardabweichung

Die Standardabweichung ist die Wurzel der Varianz. Gefühlt ist die Standardabweichung sowas wie die durchschnittliche Abweichung (Beträge) vom Mittelwert (eben durch die Quadrierung und Rückrechnung über die Wurzel nicht ganz dasselbe).

$$s_x = \hat{\sigma}_x = \sqrt{\frac{1}{n-1} \sum_i^n (x_i - \bar{x})^2}$$

Die Standardabweichung als unbekannter Parameter wird als  $\sigma_x$  gekennzeichnet und der Kennwert als  $s_x$ .



# Standardabweichung im Vergleich



<detached>

<detached>

<detached>

<detached>

<detached>

<detached>

<detached>

<detached>

<detached>

<detached>

<detached>

<detached>

<detached>

<detached>



# Standardfehler

## Definition und Eigenschaften

Die Standardabweichung einer Stichprobenkennwerteverteilung nennt man *Standardfehler* (*standard error, SE*). Wenn es sich um eine Verteilung von *Mittelwerten* handelt *Standardfehler des Mittelwerts* ( $\sigma_{\bar{x}}$ ).

Der Standardfehler gibt an, wie gut  $\bar{x}$  den Populationsparameter  $\mu$  schätzt.

### Eigenschaften

- Der Standardfehler jedes Kennwertes nimmt mit grösser werdendem  $n$  ab 🖱️  $\sigma_{\bar{x}}$  ist umgekehrt proportional zu  $\sqrt{n}$ .
- Je stärker das gemessene Merkmal in der Population streut, desto grösser der Standardfehler 🖱️  $\sigma_{\bar{x}}$  ist proportional zu  $\sigma$ .

Die Formel:

$$\begin{aligned}\sigma_{\bar{x}} &= \sqrt{\frac{\sigma^2}{n}} \\ &= \frac{\sigma}{\sqrt{n}}\end{aligned}$$





# Effekte des Stichprobenumfangs



# Standardisierung

## Standardisierung

Standardisierung von Verteilungen und Kennwerten macht alles vergleichbar.

- $z_i = \frac{x_i - \bar{x}}{s}$ .
- Eine z-Verteilte Grösse hat immer:  $\bar{x} = 0$  und  $s = 1$ .
- Standardisierte Verteilungen und Kennwerte sind vergleichbar.



# Konfidenzintervalle für Mittelwerte

## Konfidenzintervalle

Konfidenzintervalle geben einen Wertebereich an, in dem die Parameter (GG) der Stichprobenkennwerte mit einer angebbaren Wahrscheinlichkeit liegen.

$$\text{KI: } \bar{X} \pm z_1 \cdot SE$$

$$\text{KI: } \bar{X} \pm z_1 \cdot \frac{s_x}{\sqrt{n}}$$

$$\text{KI}_{l.05} = \bar{x} - 1.96 \cdot \frac{s_x}{\sqrt{n}}$$

$$\text{KI}_{r.05} = \bar{x} + 1.96 \cdot \frac{s_x}{\sqrt{n}}$$



## 3.2 Bivariate Statistik



## 3.2.1 Kovarianz und Korrelation

$$\text{cov} = C = \frac{1}{n} \sum_i^n (x_i - \bar{x})(y_i - \bar{y})$$

$$r = \frac{\sum_i^n (x_i - \bar{x})(y_i - \bar{y})}{n \cdot s_x \cdot s_y}$$



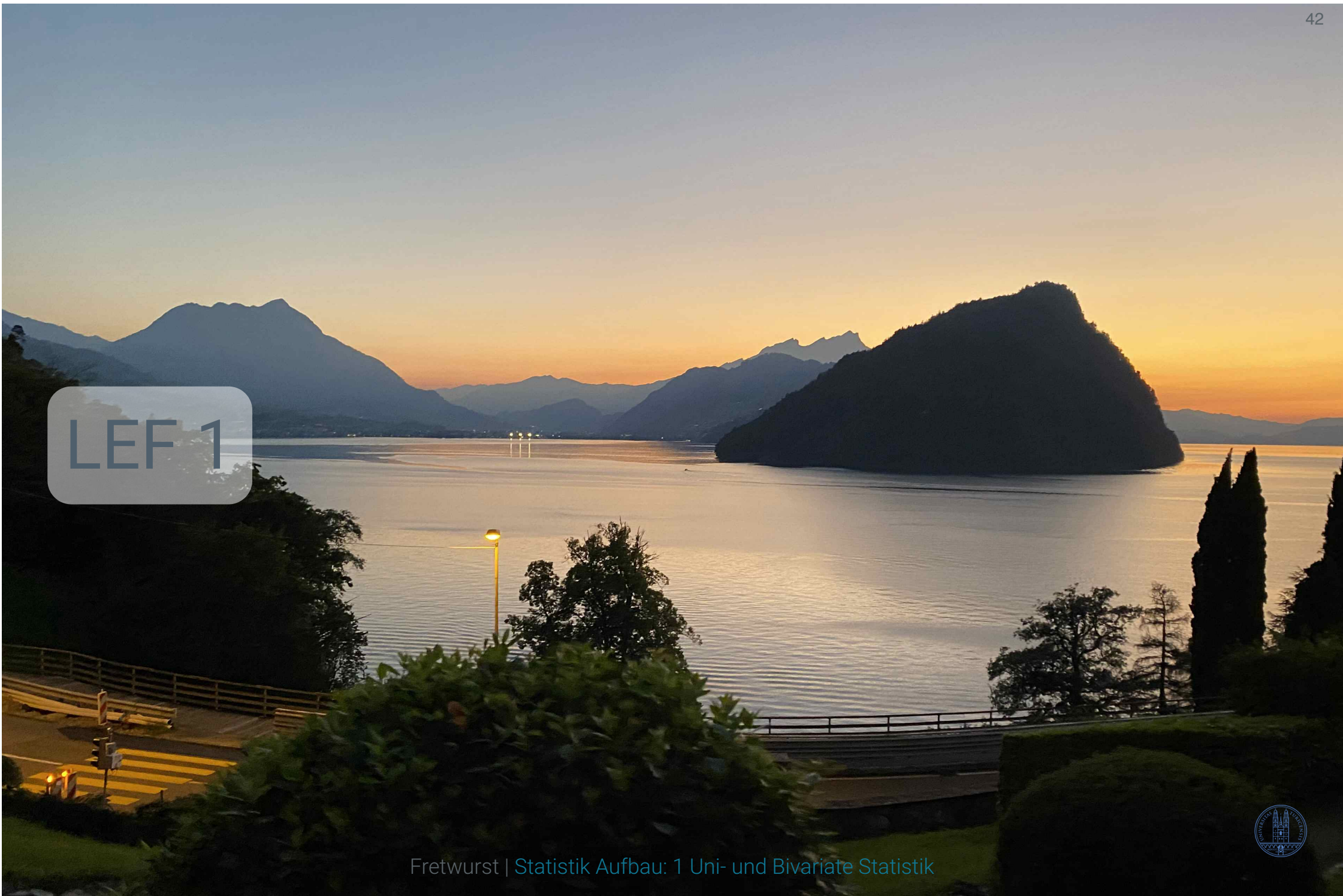
## 3.3 Hypothesentesten

### Testen gegen die Nullhypothese

1. Könnte in der Auswahlgesamtheit der wahre Wert auch 0 sein, oder ein anderes Vorzeichen haben?
2. Die Nullhypothese ist eine statistische Hypothese gegen Falschentscheidungen aufgrund von Zufallsziehungen.
3. Nullhypothesen werden anhand von bekannten Verteilungen getestet.



LEF 1



# Essayfragen

1. Was ist der Unterschied zwischen unstandardisierten und standardisierten Kennwerten?
2. Welche Masse der zentralen Tendenz kennen Sie?
3. Welche Streumasse kennen Sie?
4. Was kommt raus, wenn man die Kovarianz einer Variablen mit sich selbst berechnet?
5.
  - a. Welche Skalenniveaus kennen Sie?
  - b. Was macht eine Nominalskala aus?
  - c. Was macht eine metrische Skala aus?





# MC-Fragen



## MC 1.1.

## MC 1.1: Sind folgende Aussagen richtig oder falsch?

richtig	falsch	Aussagen
<input type="radio"/>	<input type="radio"/>	Der Mittelwert ist dasselbe wie der Durchschnitt.
<input type="radio"/>	<input type="radio"/>	Der Mittelwert einer dichotomen Variablen entspricht dem Prozentsatz der 1-Werte.
<input type="radio"/>	<input type="radio"/>	Der Mittelwert wird auch als "Mittel" oder "Arithmetisches Mittel" bezeichnet.
<input type="radio"/>	<input type="radio"/>	Je grösser ein Mittelwert ist, desto eher ist er signifikant.

Punkte: 0



## MC 1.2.

## MC 1.2: Sind folgende Aussagen richtig oder falsch?

richtig	falsch	Aussagen
<input type="radio"/>	<input type="radio"/>	Die Standardabweichung ist die standardisierte Form der Varianz.
<input type="radio"/>	<input type="radio"/>	Die Standardabweichung hat $n-1$ Freiheitsgrade.
<input type="radio"/>	<input type="radio"/>	Die Standardabweichung "s" liegt immer zwischen 0 und 1.
<input type="radio"/>	<input type="radio"/>	Die Standardabweichung ist die durchschnittliche Abweichung vom Mittelwert.

**Punkte: 0**

## MC 1.3.

**MC 1.3: Sind folgende Aussagen richtig oder falsch?**

<b>richtig</b>	<b>falsch</b>	<b>Aussagen</b>
<input type="radio"/>	<input type="radio"/>	Die Covarianz ist das Quadrat der Korrelation.
<input type="radio"/>	<input type="radio"/>	Die Covarianz ist skalenabhängig und kann daher negativ oder positiv und unendlich klein oder gross sein.
<input type="radio"/>	<input type="radio"/>	Die Korrelation "r" liegt immer zwischen 0 und 1.
<input type="radio"/>	<input type="radio"/>	Eine Korrelation von genau 0 kann nie signifikant sein.

**Punkte: 0**

## MC 1.4.

## MC 1.4: Sind folgende Aussagen richtig oder falsch?

- | richtig               | falsch                | Aussagen   |
|-----------------------|-----------------------|--|
| <input type="radio"/> | <input type="radio"/> | Bei der bivariaten Korrelation sind $r$ und die Wurzel aus $R^2$ identisch.  |
| <input type="radio"/> | <input type="radio"/> | Das $b$ einer bivariaten Regressionsgeraden liegt immer zwischen -1 und 1.   |
| <input type="radio"/> | <input type="radio"/> | Wenn das $b$ nicht signifikant ist, kann nicht mit 95% Wahrscheinlichkeit ausgeschlossen werden, dass es in Wirklichkeit Null ist oder ein zum Stichprobenwert entgegengesetztes Vorzeichen hat. |
| <input type="radio"/> | <input type="radio"/> | Wenn ein $b$ signifikant ist, dann ist auch BETA signifikant.  |

Punkte: 0



## MC 1.5.

**MC 1.5: Sind folgende Aussagen richtig oder falsch?**

<b>richtig</b>	<b>falsch</b>	<b>Aussagen</b>
<input type="radio"/>	<input type="radio"/>	Statistische Signifikanz bedeutet im Grunde wissenschaftliche Relevanz.
<input type="radio"/>	<input type="radio"/>	Bei kleinen Stichproben können Ergebnisse auch schon mal signifikant werden, obwohl die Effekte so kleine sind, dass sie zu vernachlässigen sind.
<input type="radio"/>	<input type="radio"/>	Je grösser $n$ , desto schneller wird derselbe Effekt (z.B. Mittelwertunterschied) signifikant.
<input type="radio"/>	<input type="radio"/>	Wenn man einmal ein Signifikanzniveau (z.B. 95%) festgelegt hat, sollte man auch dabei bleiben.

**Punkte: 0**

## MC 1.6.

**MC 1.6: Sind folgende Aussagen richtig oder falsch?**

<b>richtig</b>	<b>falsch</b>	<b>Aussagen</b>
<input type="radio"/>	<input type="radio"/>	Wenn das Konfidenzintervall eines Kennwertes kleiner als .05 ist, dann ist es signifikant.
<input type="radio"/>	<input type="radio"/>	Wenn das Konfidenzintervall eines Mittelwertes die 0 nicht einschliesst, dann ist der Mittelwert signifikant von 0 verschieden.
<input type="radio"/>	<input type="radio"/>	Wählt man ein höheres Signifikanzniveau (z.B. 99% statt 95%), dann wird das Konfidenzintervall breiter.
<input type="radio"/>	<input type="radio"/>	Je breiter ein Konfidenzintervall, desto besser ist ein Kennwert interpretierbar.

**Punkte: 0**

**Für LEF 1: 0 von 12 Punkten, was 0% und etwa einer 1 entspricht.**



# Take Home

## Statistik

- ist ein Modell von Realität
- ist eine mächtige Denkweise
- ist Basis und Kern von Data Science

## Aus Statistik Einführung muss sitzen

- Mittelwerte
- Varianz/ Standardabweichung  $s$ - Covarianz/ Korrelation





# Ausblick

## Grundlagen der Modellbildung

Regression mit zwei Unabhängigen

## Inferenzstatistische Grundlagen

Punktschätzung, Intervallschätzung, Wahrscheinlichkeitsverteilungen

## Grundannahmen von OLS-Schätzungen

mässige Multikollinearität, Homoskedastizität, Modellspezifikation, Linearität



