

Statistik für Data Scientists

Vorlesung 7: Hypothesentests

Prof. Dr. Siegfried Handschuh
DS-NLP | Universität St. Gallen

Recap & Ziele heute

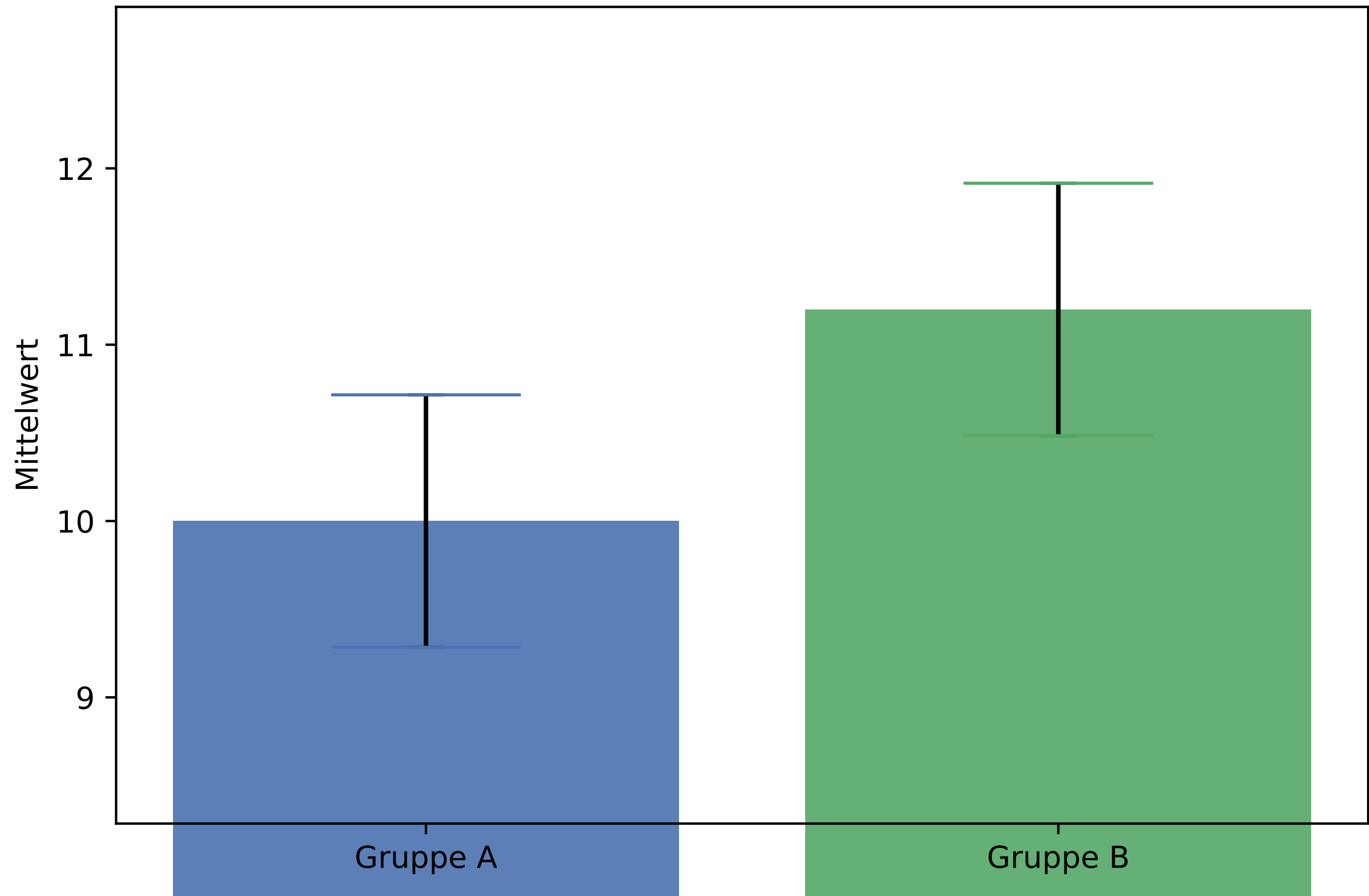
- Rückblick VL6: Schätzen & Konfidenzintervalle.
- Heute: Grundlagen der Hypothesentests.
- Null- und Alternativhypothese (H_0, H_1).
- Fehlerarten & Teststärke (Power).
- p-Wert richtig interpretieren.
- t-Test an Beispielen.
- Kritik & gute Praxis.

Motivation & Grundidee

Warum Hypothesentests?

- Frage: Ist ein Unterschied **zufällig** oder **echt**?
- Beispiele:
 - Wirkt ein Medikament tatsächlich?
 - Funktioniert A/B-Variante besser?
 - Verdienen Männer und Frauen gleich viel?
- Statistik liefert ein objektives Verfahren, um Unsicherheit zu quantifizieren.

Motivation: Unterschied oder Zufall?

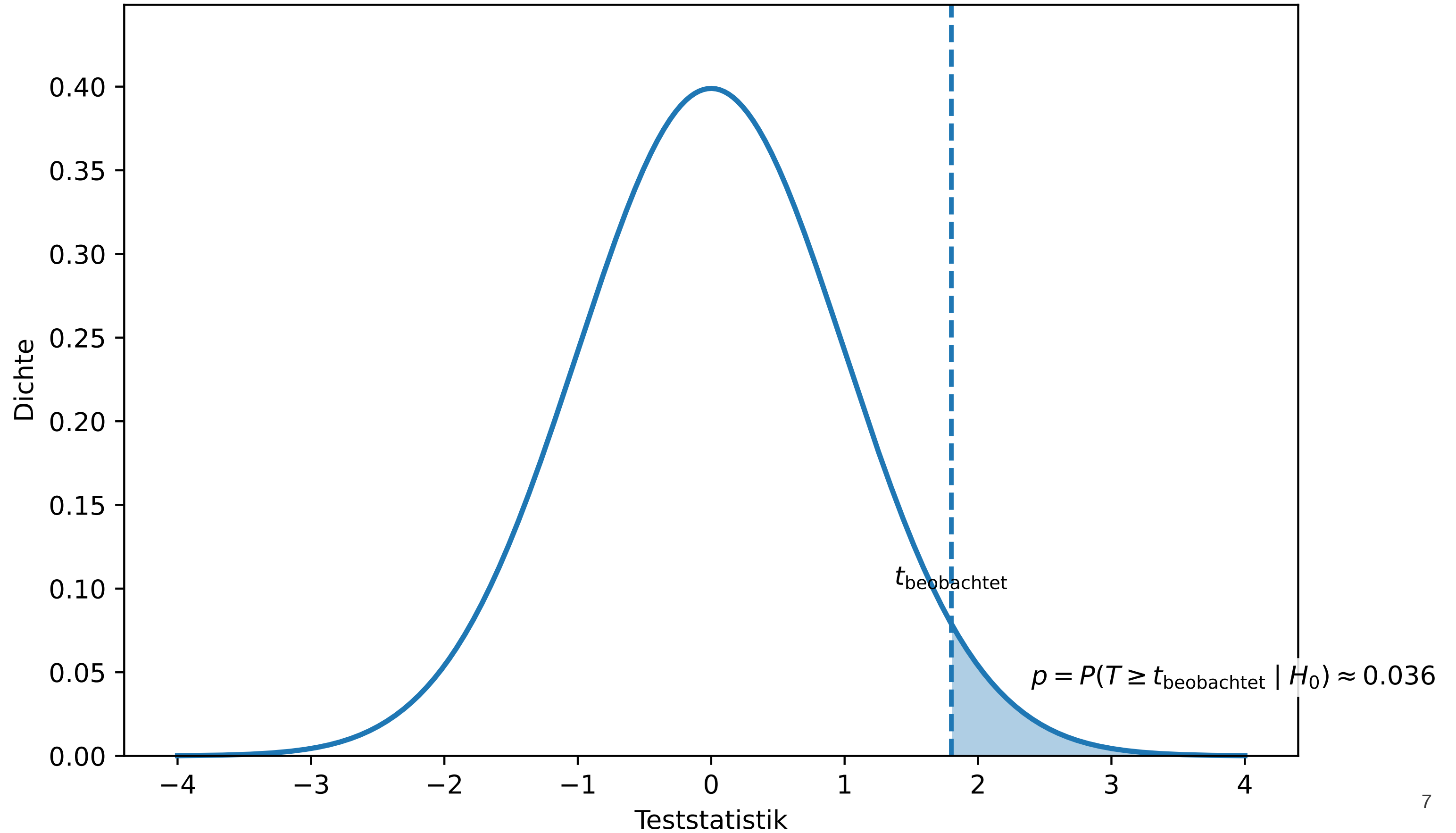


Grundidee des Tests

- Formuliere H_0 : kein Effekt.
- Formuliere H_1 : Effekt vorhanden.
- Wähle eine Teststatistik T .
- Vergleiche den beobachteten Wert $t_{\text{beobachtet}}$ mit der Referenzverteilung unter H_0 .
- Kleine p-Werte \rightarrow starke Evidenz gegen H_0 .

$$p = P(T \geq t_{\text{beobachtet}} \mid H_0)$$

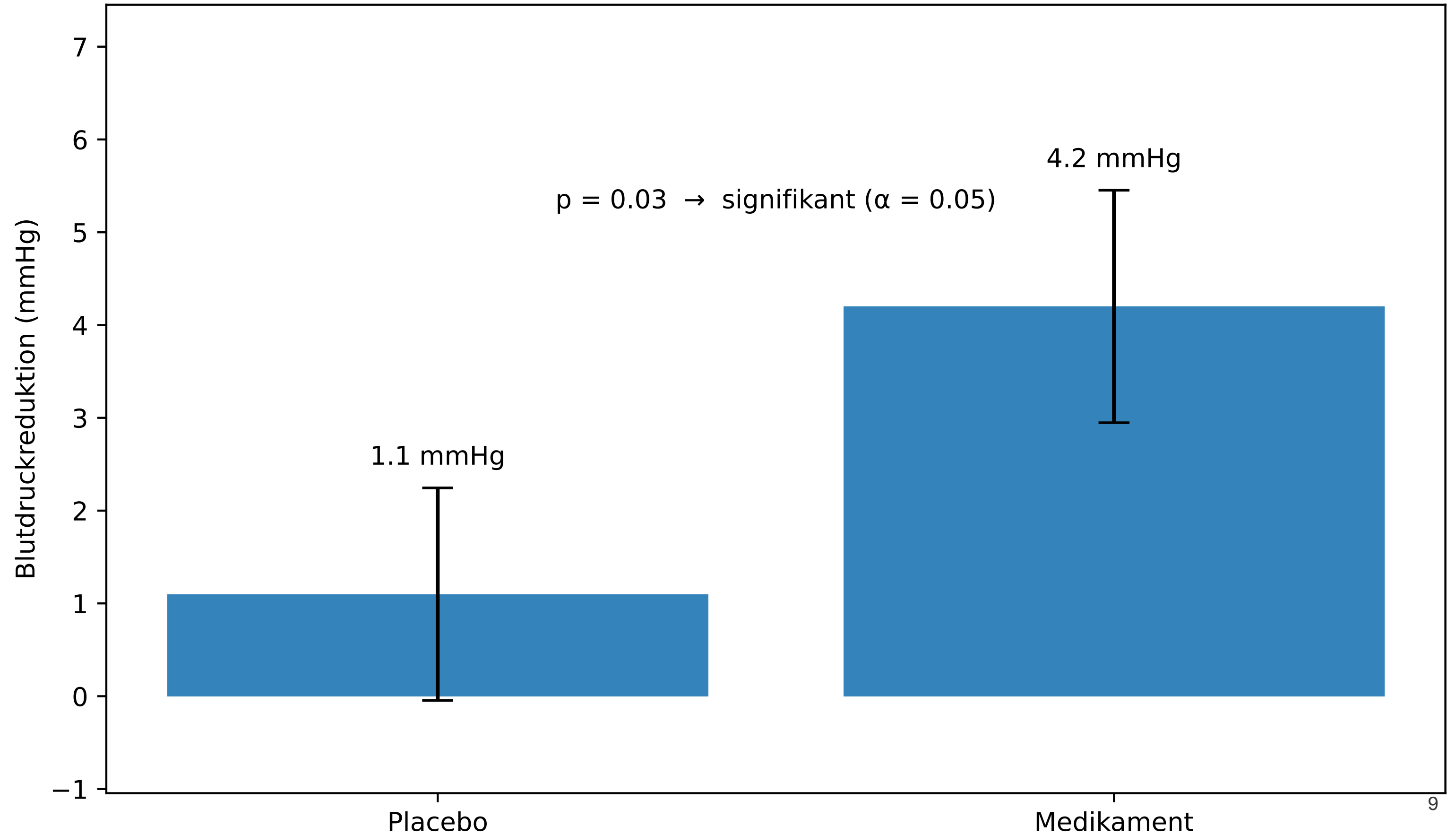
p-Wert unter H_0 (rechte Flanke)



Beispiel: Medikamententest

- Zwei Gruppen: Medikament vs. Placebo.
- Mittelwert der Blutdruckreduktion: 4.2 mmHg vs. 1.1 mmHg.
- t-Test $\rightarrow p = 0.03$ bei $\alpha = 0.05 \rightarrow H_0$ verwerfen.
- Signifikant? Ja; praktisch relevant? Fraglich.

Beispiel: Medikamententest



Take-away: Motivation & Grundidee

- Hypothesentests = Werkzeug gegen Zufallsurteile.
- Formale Bausteine: H_0 , H_1 , Teststatistik T , p-Wert.
- p-Wert misst Seltenheit der Daten unter H_0 , nicht deren Wahrheit.
- Beispiel: Medikament – signifikant, aber nicht immer relevant.

Statistik liefert Entscheidungen unter Unsicherheit: nicht absolute Wahrheiten.

Fehlerarten & Power

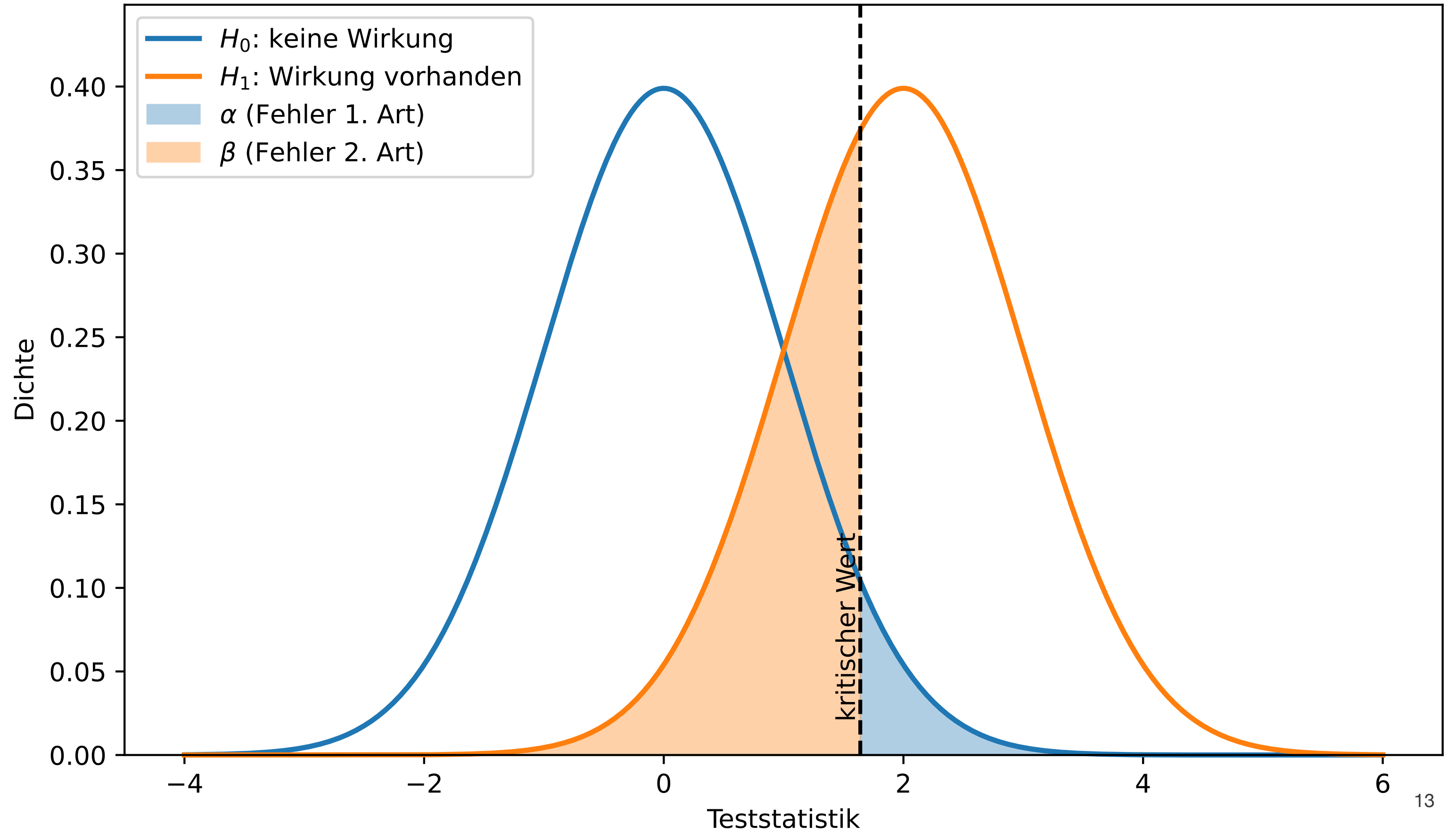
Fehlerarten

Typ	Beschreibung	Symbol
Fehler 1. Art	H_0 ist wahr, wird aber verworfen. <i>⟨False Positive⟩ → Fehlalarm.</i>	α
Fehler 2. Art	H_0 ist falsch, wird aber nicht verworfen. <i>⟨False Negative⟩ → übersehener Effekt.</i>	β

Teststärke (Power): $1 - \beta$

Ziel: α klein halten, Power gross halten.

Fehlerarten bei Hypothesentests



Teststärke (Power)

Definition: Wahrscheinlichkeit, einen echten Effekt zu entdecken

$$\text{Power} = 1 - \beta$$

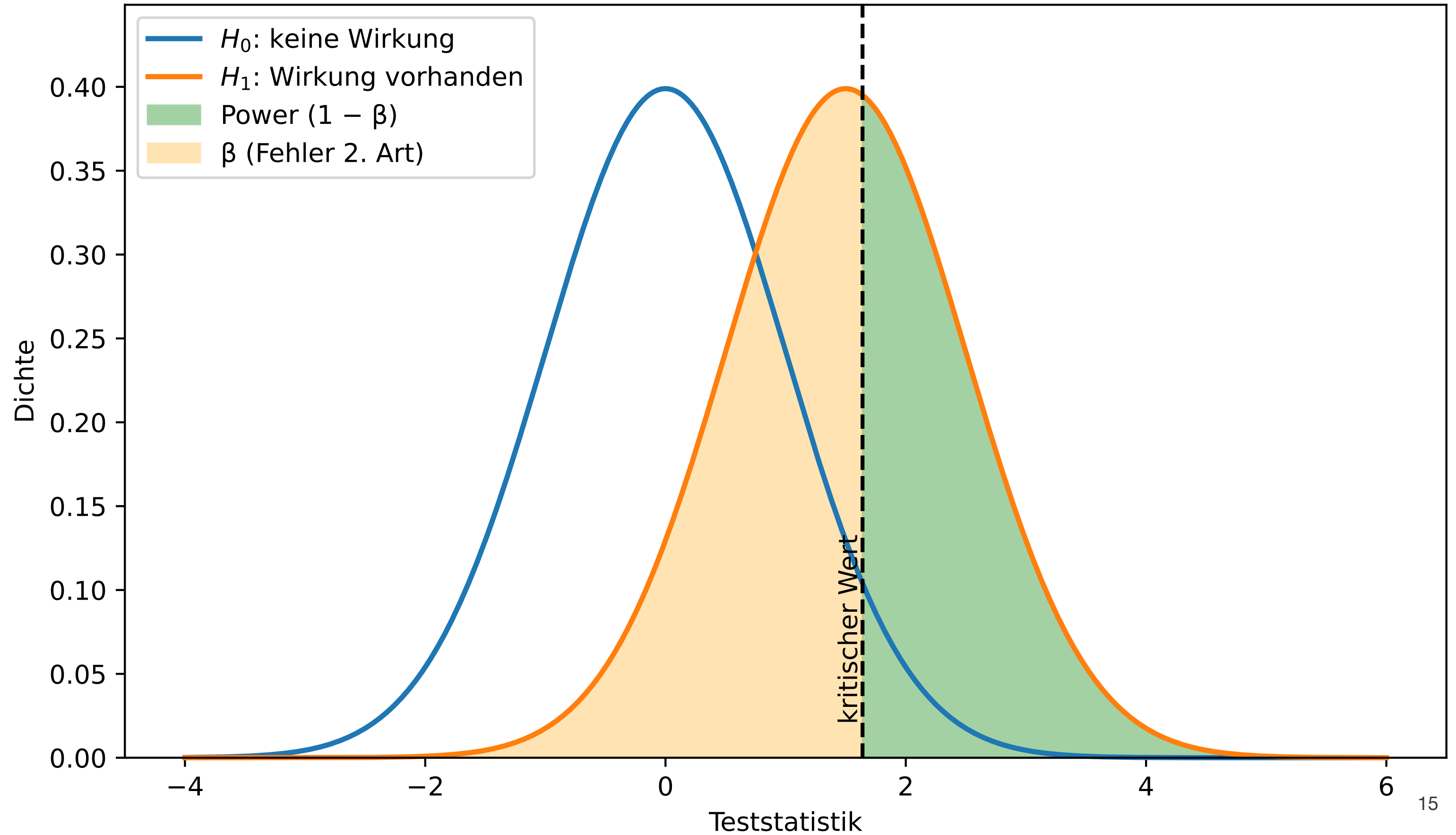
Beispiel:

Ein Medikament wirkt tatsächlich.

- Power 80 % → In 8 von 10 Studien wird die Wirkung erkannt.
- Power 20 % → In 8 von 10 Studien wird sie übersehen.

Faustregel: Power ≥ 0.80 anstreben

Teststärke (Power)



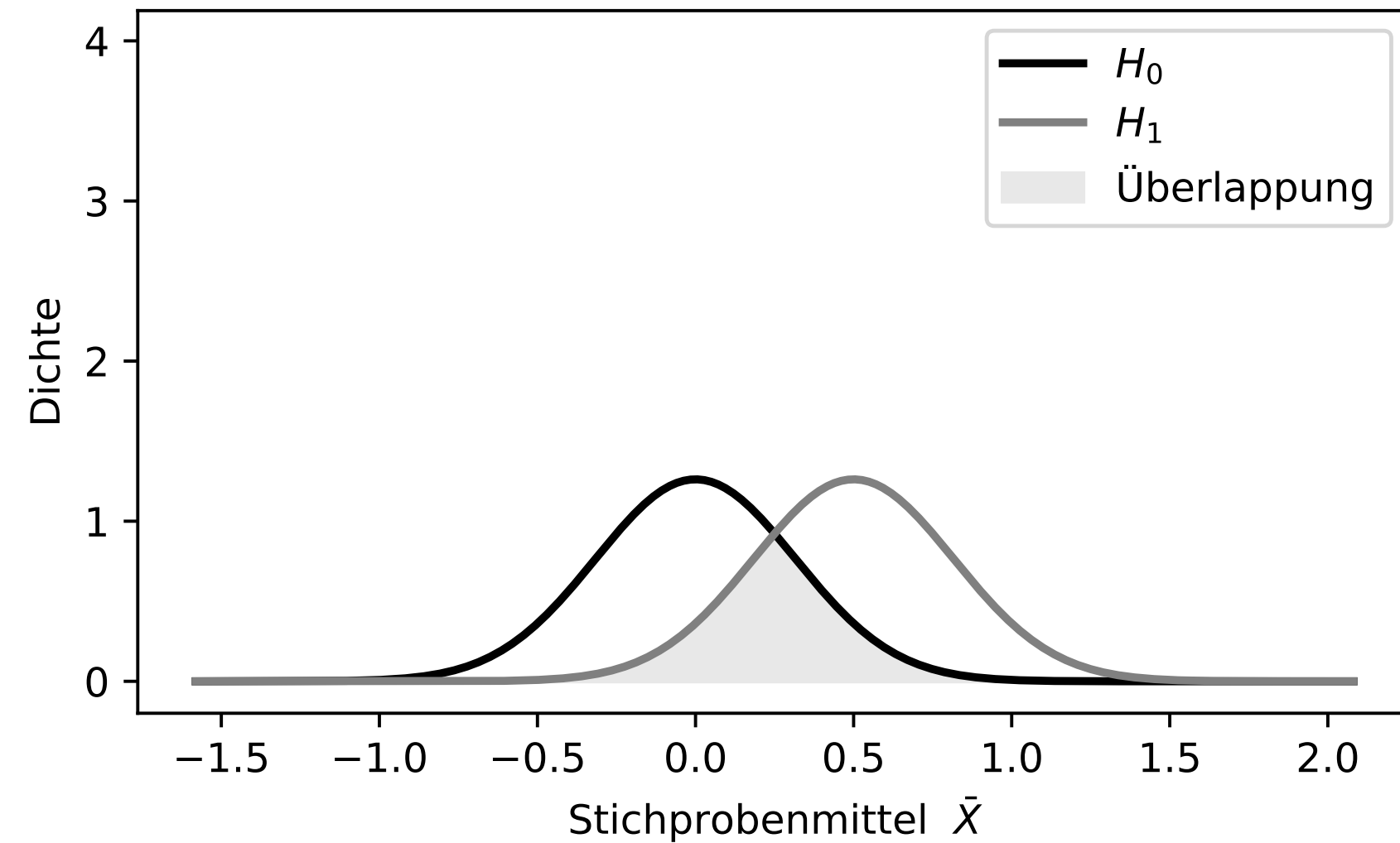
Einfluss auf Power

- **Stichprobengrösse** n : grösser \rightarrow kleinerer Standardfehler \rightarrow höhere Power.
- **Effektgrösse** Δ : grösser \rightarrow leichter erkennbar.
- **Signifikanzniveau** α : höher \rightarrow mehr Power, aber auch mehr Fehlalarme.

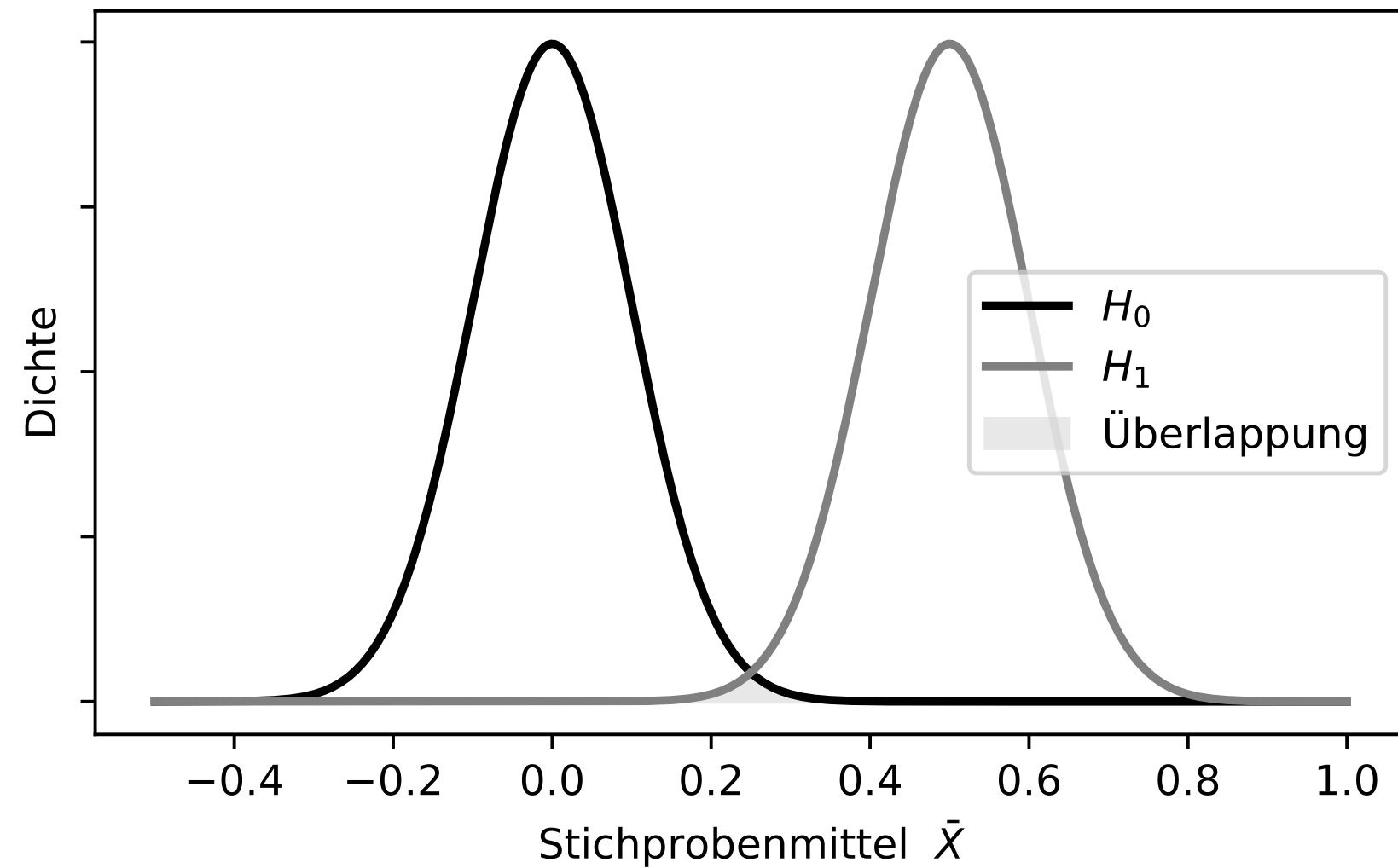
$$\text{Power} = f(n, \alpha, \Delta, \sigma)$$

Grössere Stichprobe → schmalere Kurven → geringere Überlappung → höhere Power

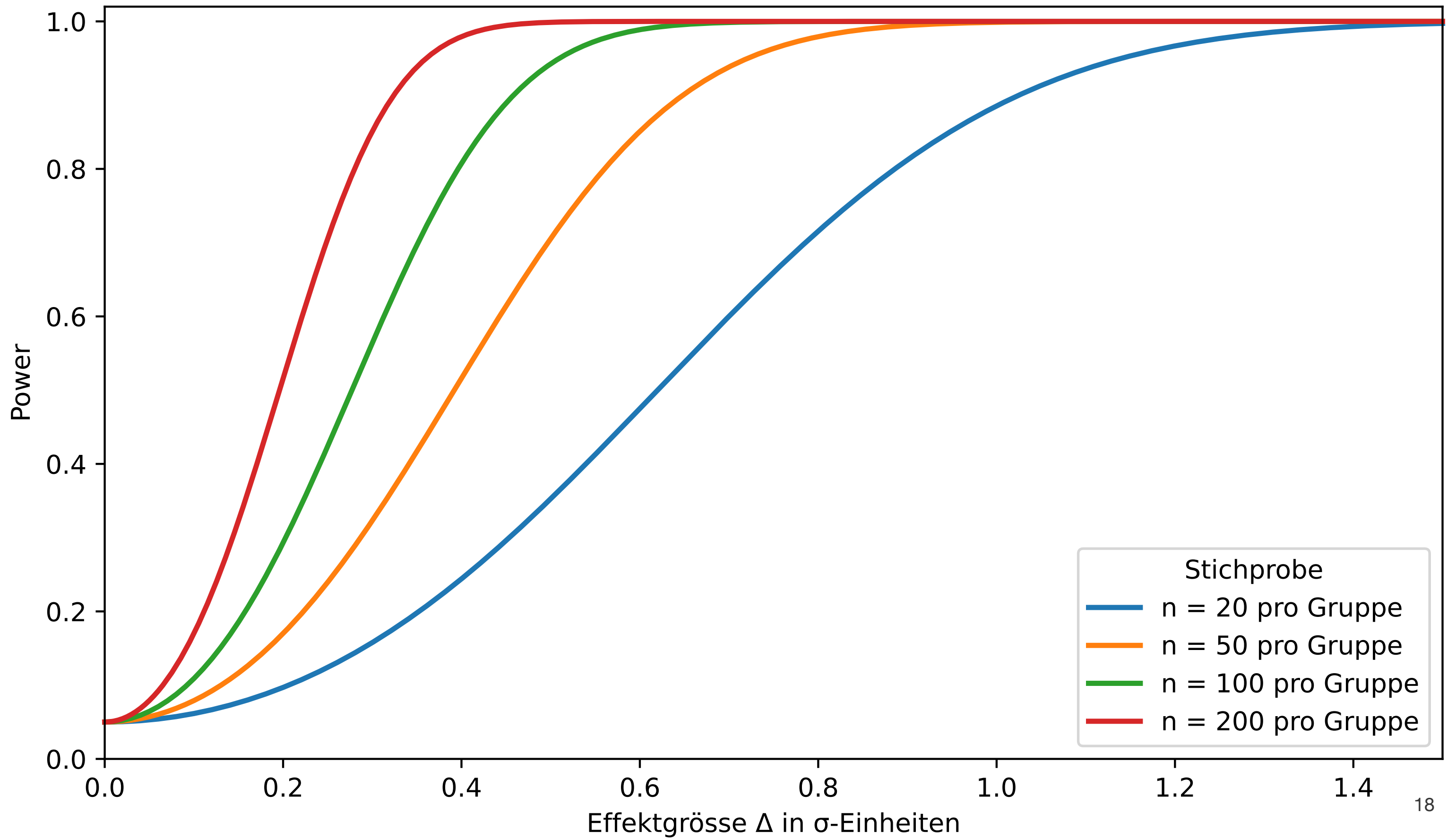
$n = 10$
 $SE = 0.316$



$n = 100$
 $SE = 0.100$



Einfluss auf Power ($\alpha = 0.05$)



Take-away: Fehlerarten & Power

- Fehler 1 (Typ I): H_0 wahr, aber verworfen $\rightarrow \alpha$ (Fehlalarm).
- Fehler 2 (Typ II): H_0 falsch, aber beibehalten $\rightarrow \beta$ (übersehen).
- Power = $1 - \beta$: Wahrscheinlichkeit, echten Effekt zu entdecken.
- Einflussfaktoren auf Power: $n \uparrow$, Effektgrösse $\Delta \uparrow$, $\alpha \uparrow$
(Achtung: $\alpha \uparrow \Rightarrow$ mehr Fehlalarme).

Take-away:

α und β sind keine Fehler, sondern Designentscheidungen \rightarrow gute Statistik balanciert Risiko und Sensitivität.

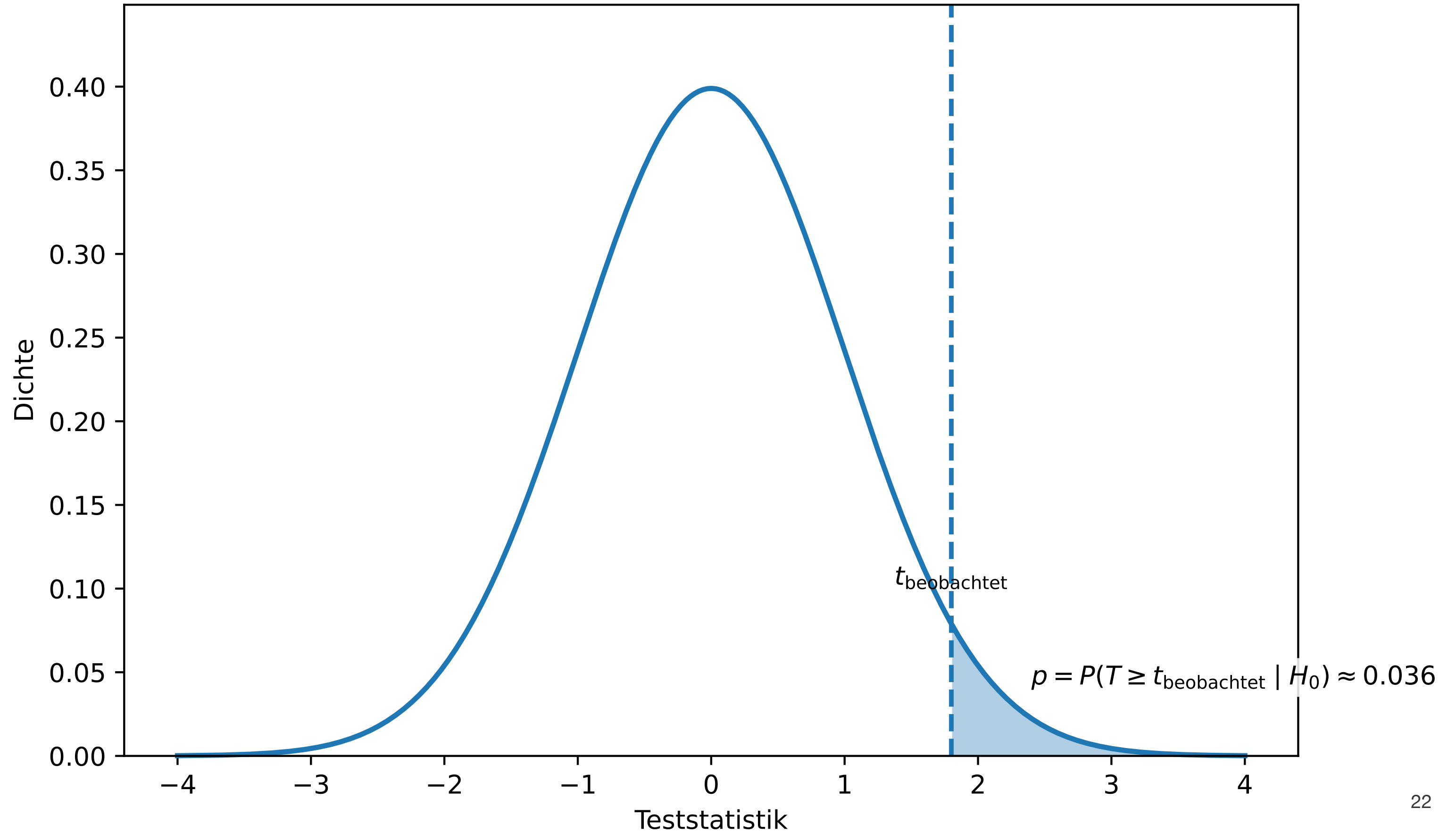
p-wert

p-Wert verstehen

- Der **p-Wert** ist die Wahrscheinlichkeit, unter H_0 ein Ergebnis **so extrem oder extremer** zu beobachten.
- Er misst die **Kompatibilität der Daten mit H_0** , nicht die Wahrscheinlichkeit, dass H_0 wahr ist.
- Kleine $p \rightarrow$ Daten **ungewöhnlich** unter H_0 .
- Grosse $p \rightarrow$ Daten **typisch** unter H_0 .

$$p = P(T \geq t_{\text{beobachtet}} \mid H_0)$$

p-Wert unter H_0 (rechte Flanke)



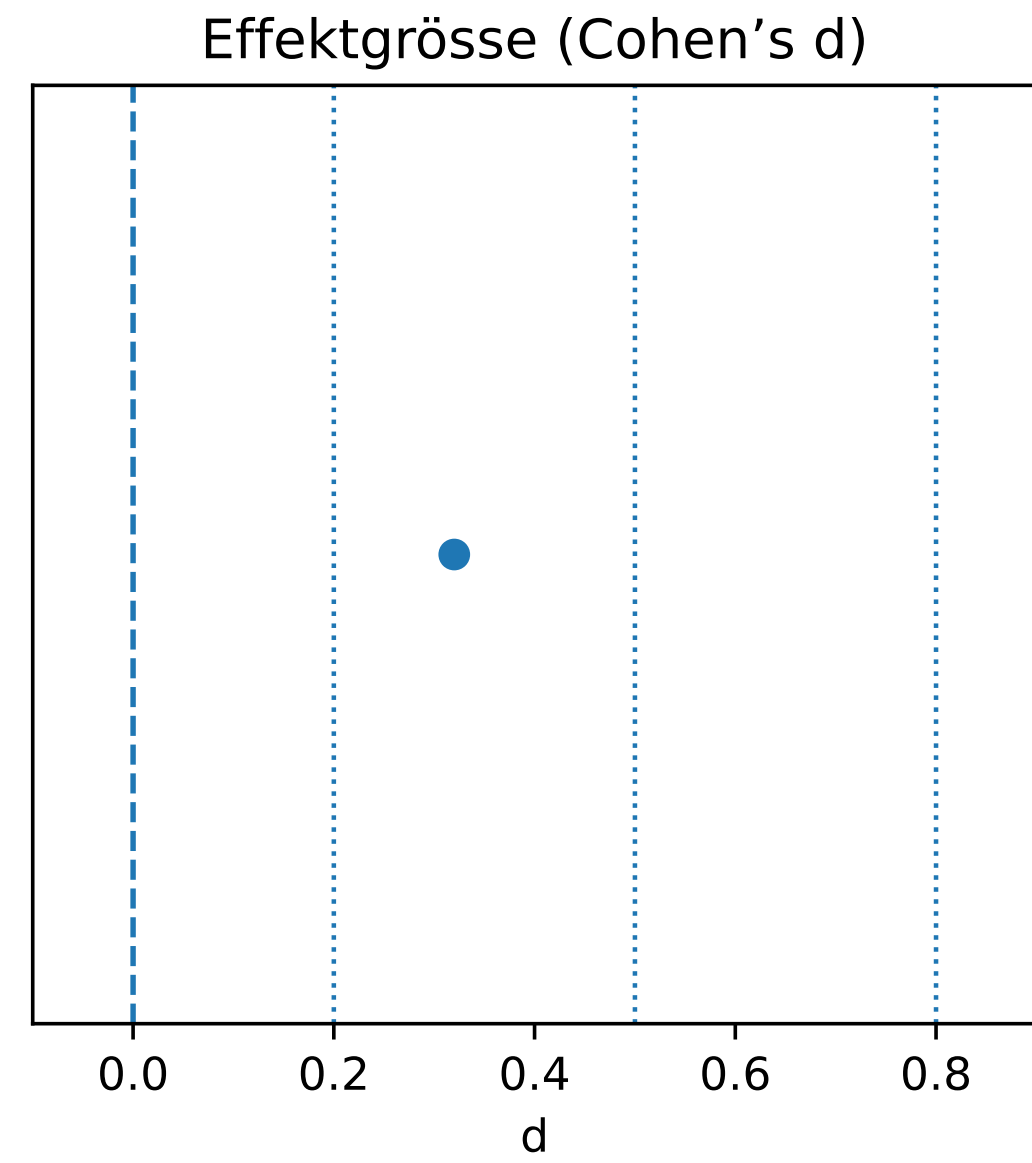
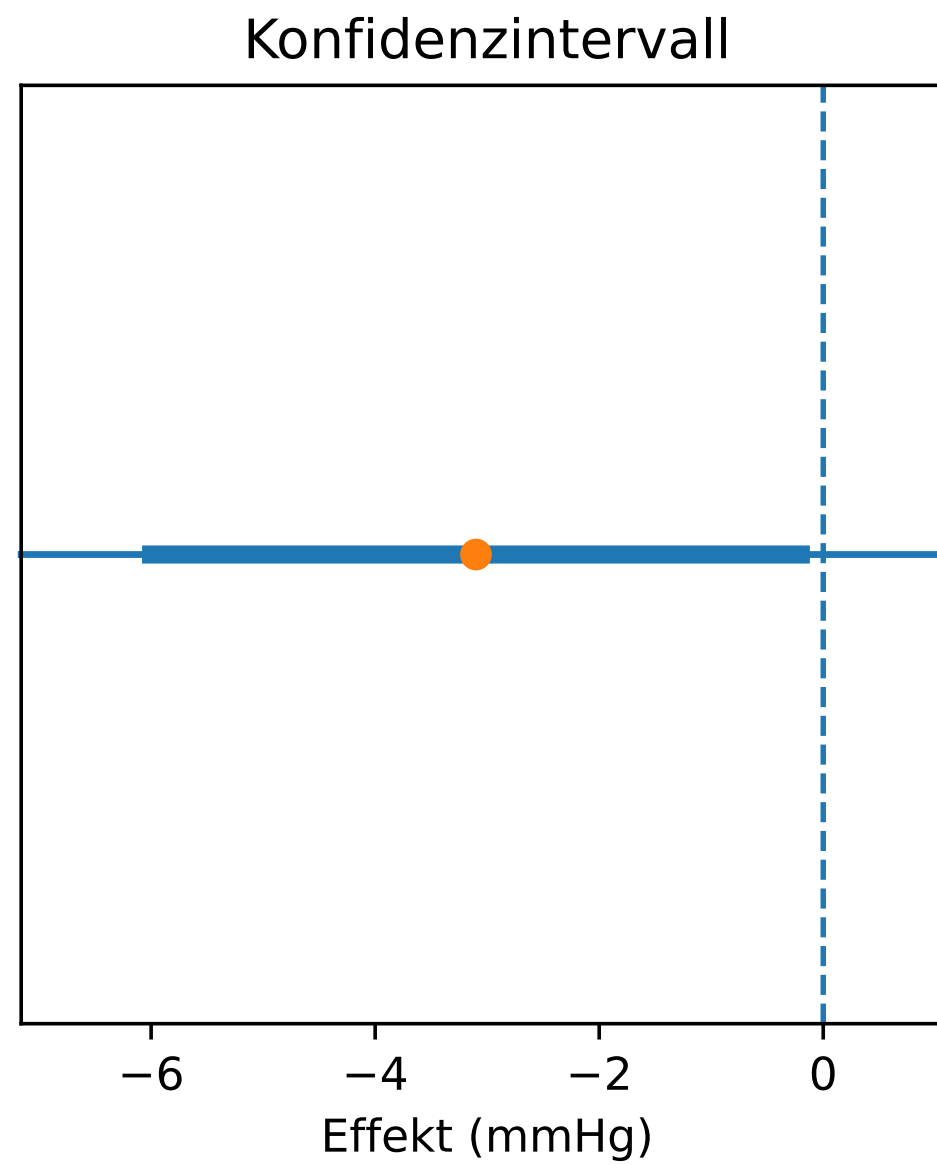
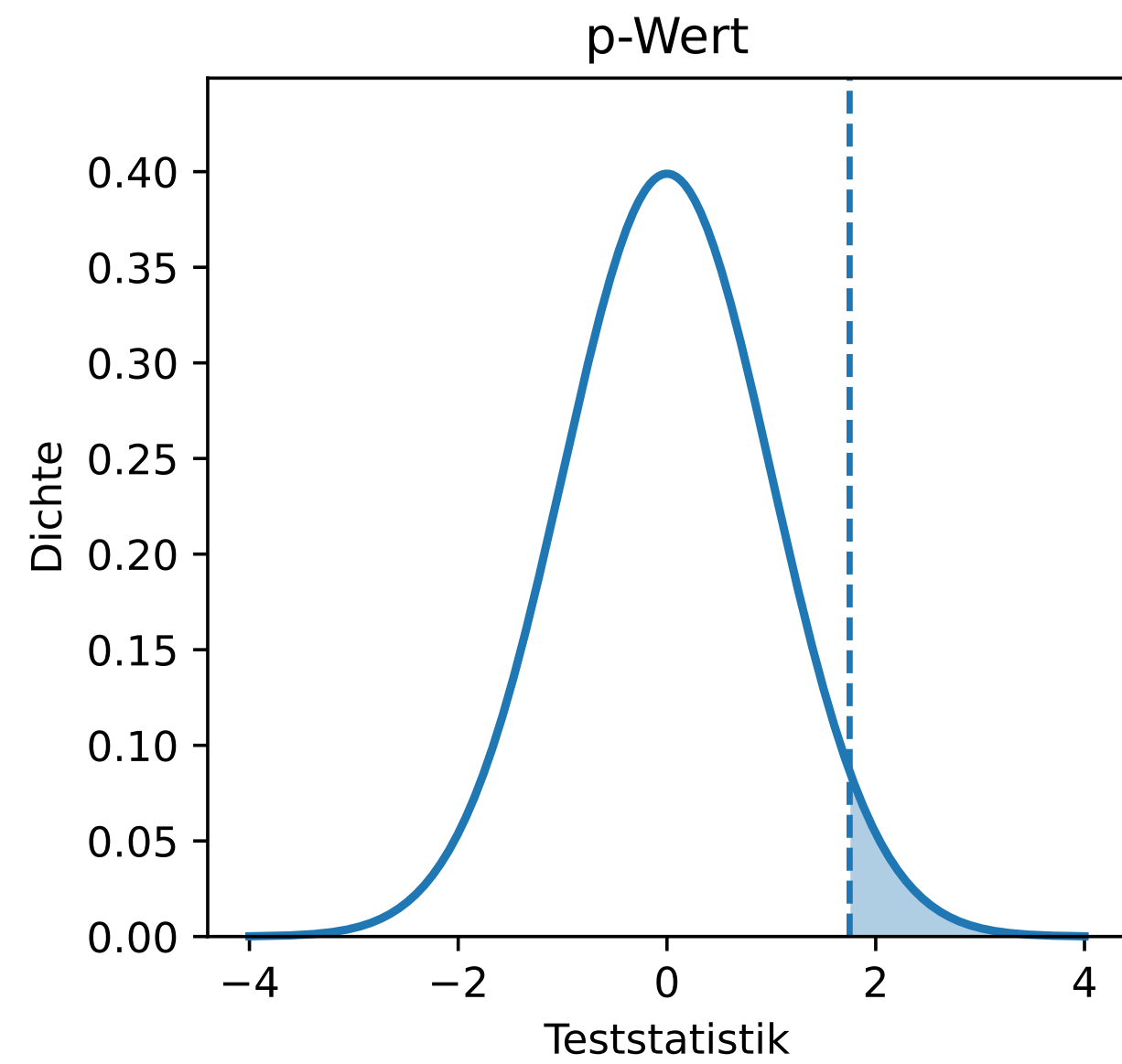
Typische Missverständnisse

- **✗** $p = 0.03 \rightarrow$ „ H_0 ist zu 97 % falsch.“
- **✗** $p = 0.20 \rightarrow$ „ H_0 stimmt.“
- **✓ Korrekt:** Der p-Wert beschreibt, wie ungewöhnlich die Daten unter H_0 sind.
- Je kleiner p , desto weniger plausibel sind die Daten, wenn H_0 wahr wäre.

Beispiel: Medikamentenstudie

- Studie: Medikament senkt Blutdruck.
- Ergebnis: $p = 0.04 \rightarrow H_0$ verwerfen (bei $\alpha = 0.05$).
- Aber: **Effektgrösse** und **Konfidenzintervall** wichtiger als p allein.
- Immer p , **CI** und **Effektgrösse** gemeinsam berichten.

Cohen's d: Effektgrösse \rightarrow **StatCoach fragen**



Take-away: p-Wert

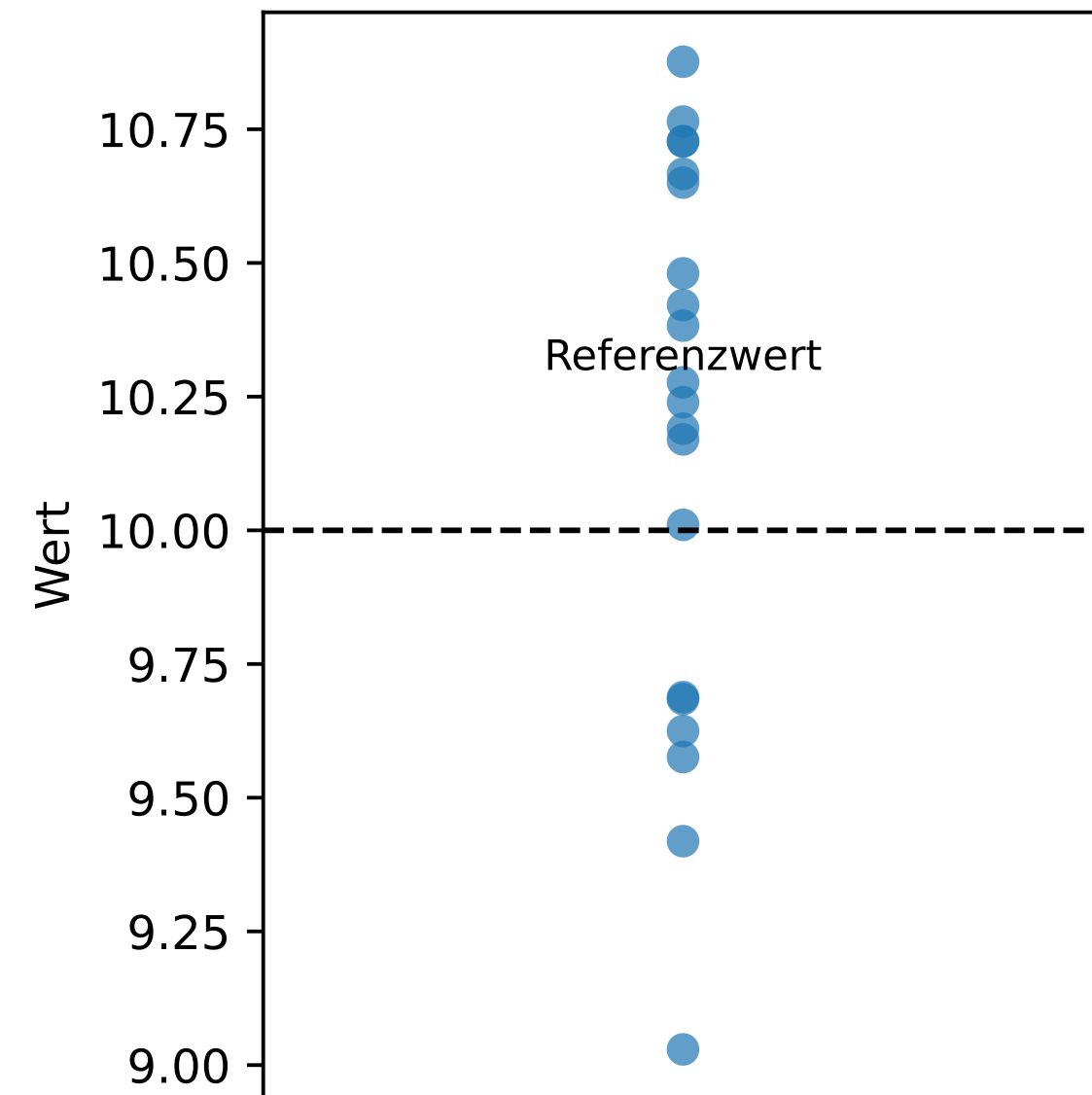
- **p-Wert:** $P(T \geq t_{\text{obs}} \mid H_0)$: **Seltenheit der Daten unter H_0 .**
- **Kein** Maß dafür, dass H_0 falsch/wahr ist (kein posterior).
- **Kleine p** \rightarrow Daten sind unter H_0 **ungewöhnlich**;
grosse p \rightarrow **typisch** unter H_0 .
- Berichte immer **p + Konfidenzintervall + Effektgrösse** (z. B. d):
Kontext entscheidet.

t-Tests

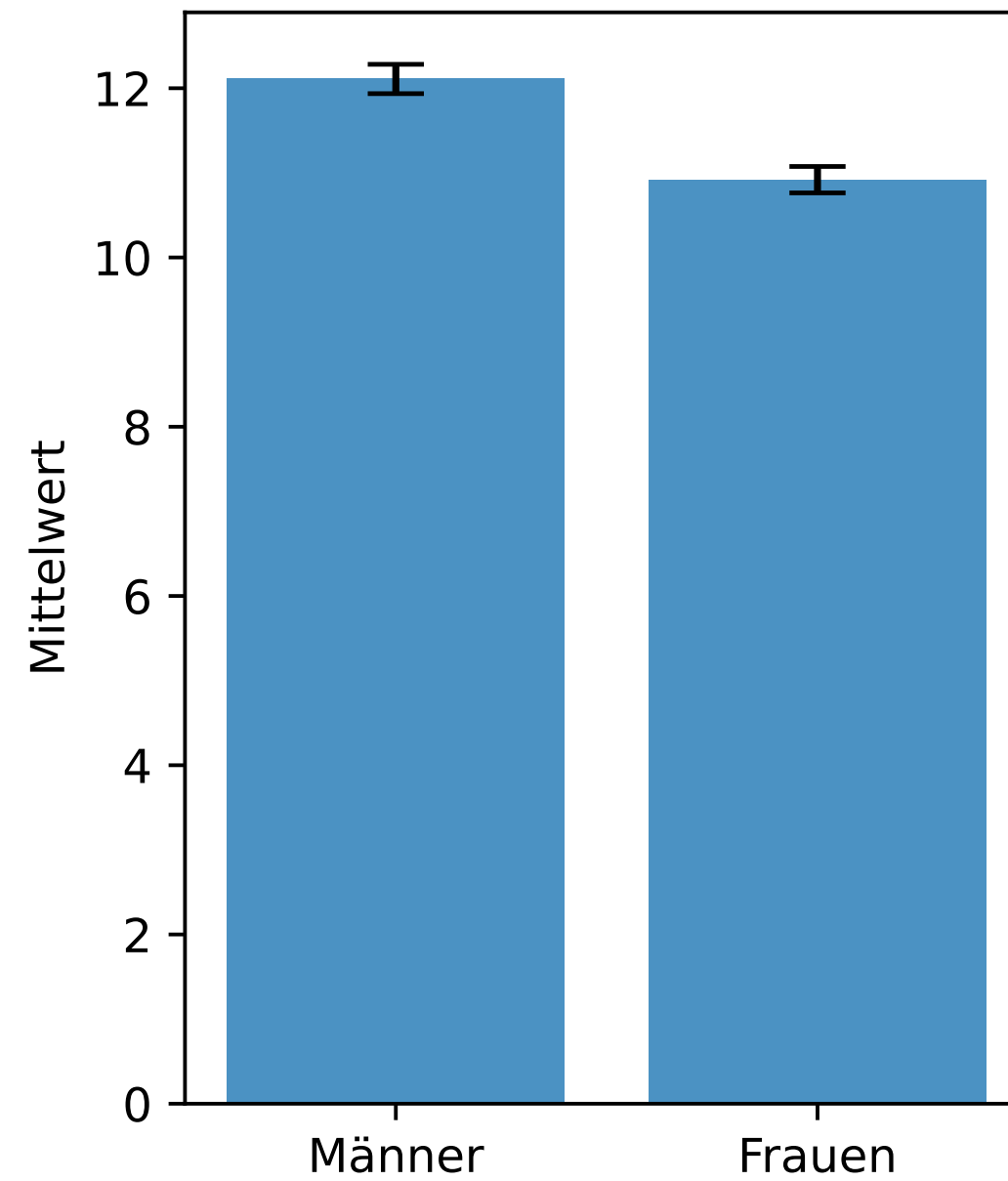
Überblick t-Tests

Testtyp	Anwendung	Beispiel
1-Sample	Mittelwert einer Stichprobe $\mu = 10$? vs. Referenzwert	
2-Sample	Vergleich zweier unabhängiger Gruppen	Männer vs. Frauen
Paired	Vergleich innerhalb von Personen (Vorher– Nachher)	Trainingseffekt

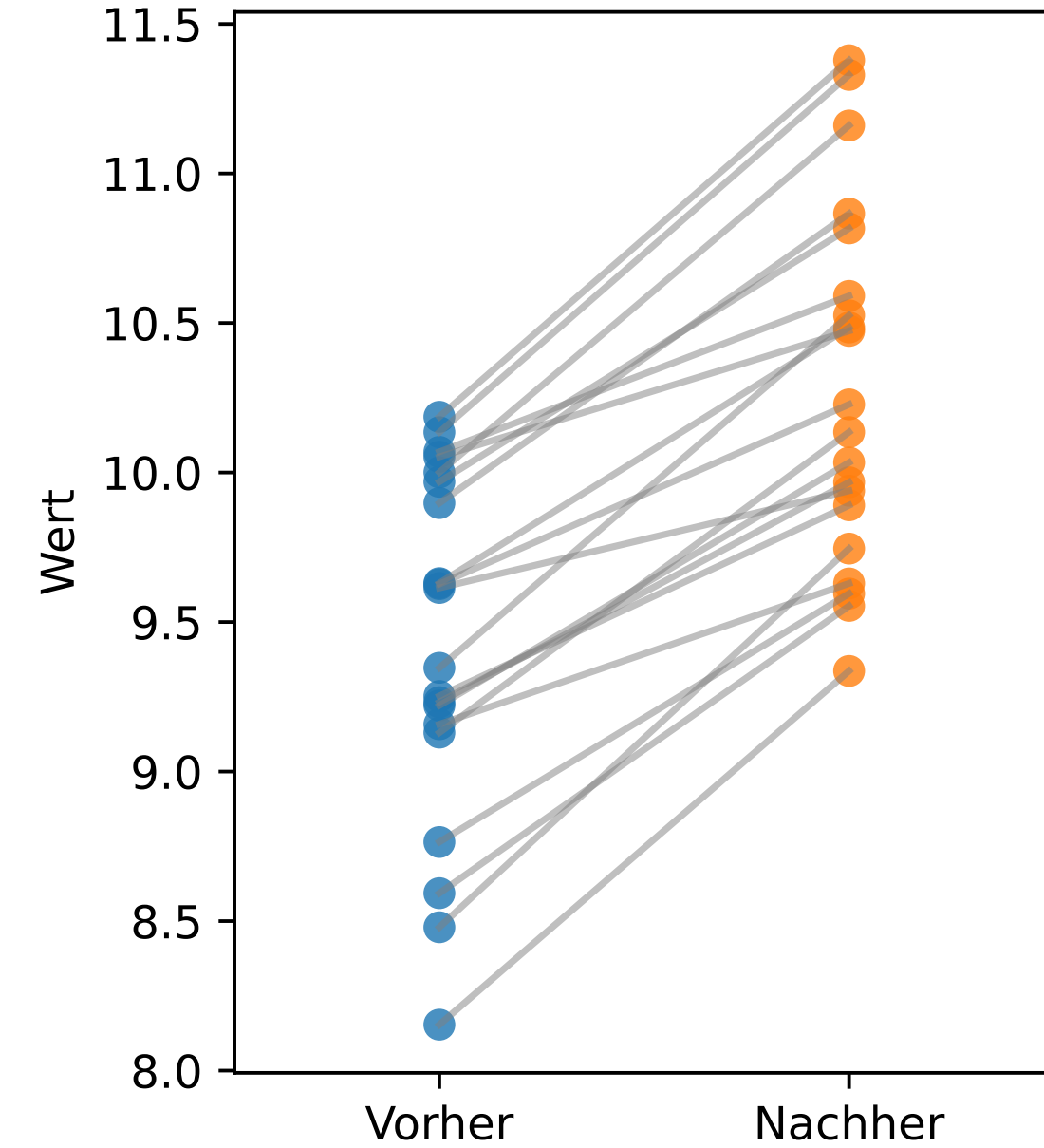
1-Sample



2-Sample



Paired



Beispiel: Zwei-Stichproben-t-Test

- Prüft, ob sich zwei Gruppen **im Mittelwert** unterscheiden.
- $H_0: \mu_1 = \mu_2$
- $H_1: \mu_1 \neq \mu_2$

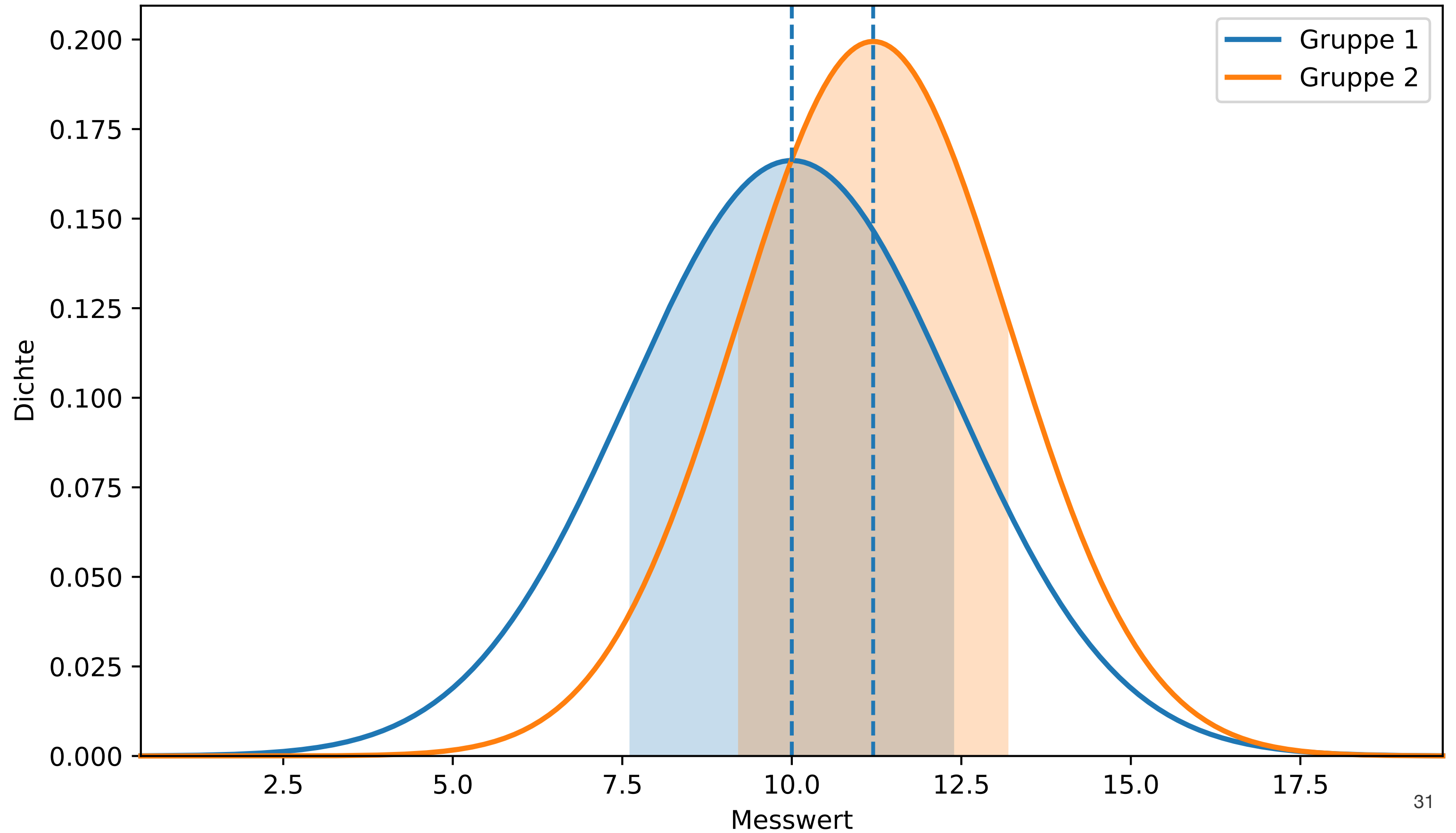
Teststatistik:

$$t = \frac{\bar{x}_1 - \bar{x}_2}{s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

mit

$$s_p^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}$$

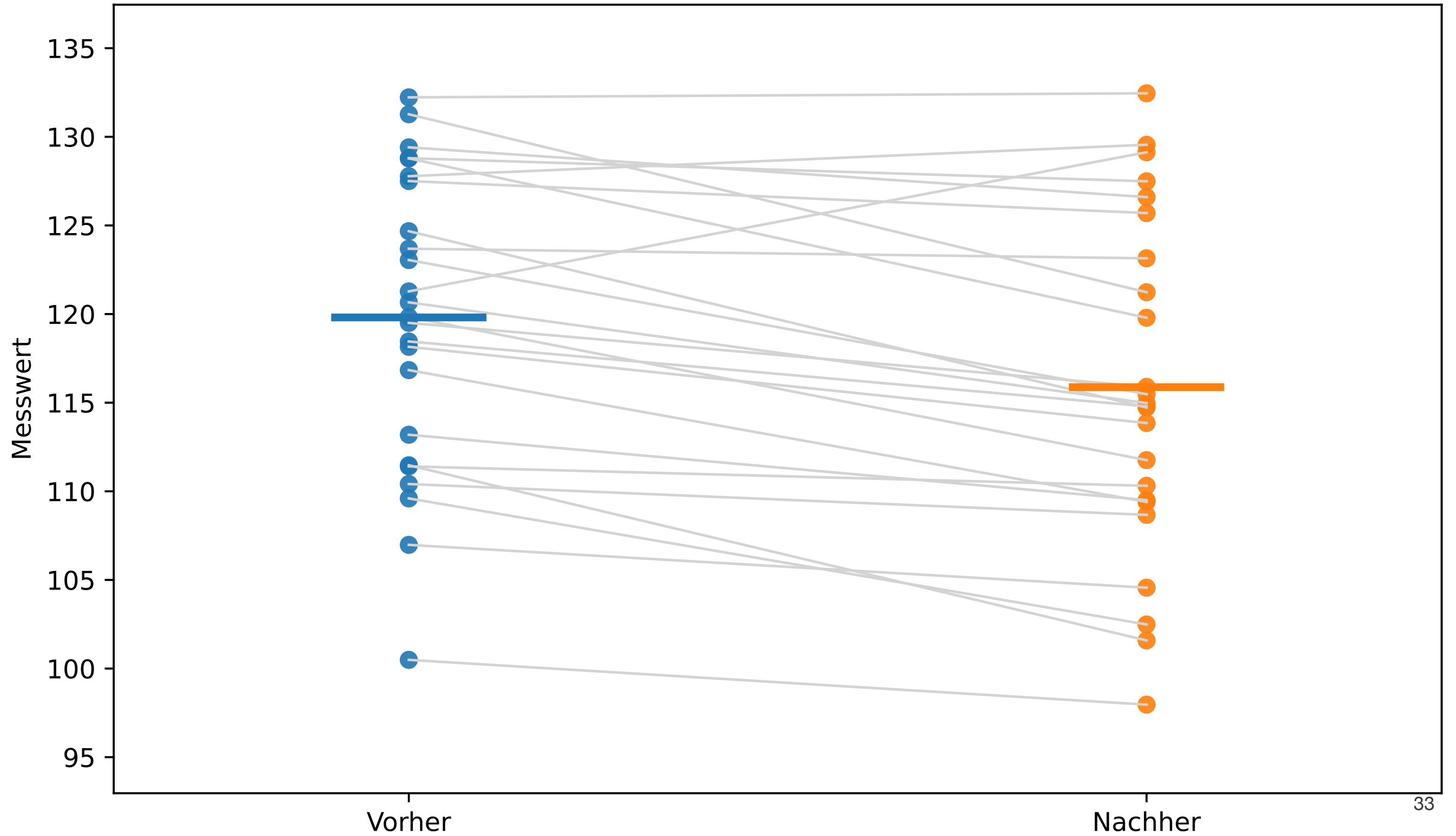
Zwei-Stichproben-t-Test: Verteilungen, Mittelwerte, ± 1 SD



Gepaarter t-Test

- Für **Vorher–Nachher-Daten** oder **gepaarte Beobachtungen**.
- Testet, ob der Mittelwert der Differenzen $\bar{d} \neq 0$ ist.
- $H_0: \mu_d = 0$
- Teststatistik: $t = \frac{\bar{d}}{s_d / \sqrt{n}}$

Gepaarter t-Test: Vorher-Nachher



Interpretation von t-Tests

Berichte immer:

- Mittelwerte und Standardabweichungen.
- Freiheitsgrade, t-Wert, p-Wert.
- Effektgrösse (z. B. Cohen's d) und Konfidenzintervall.

Beispiel-Bericht:

$t(58) = 2.13$, $p = 0.037$, $d = 0.55 \rightarrow$ mittlerer Effekt.

Ergebnis \neq Schlussfolgerung \rightarrow Kontext zählt.

Take-away: t-Tests

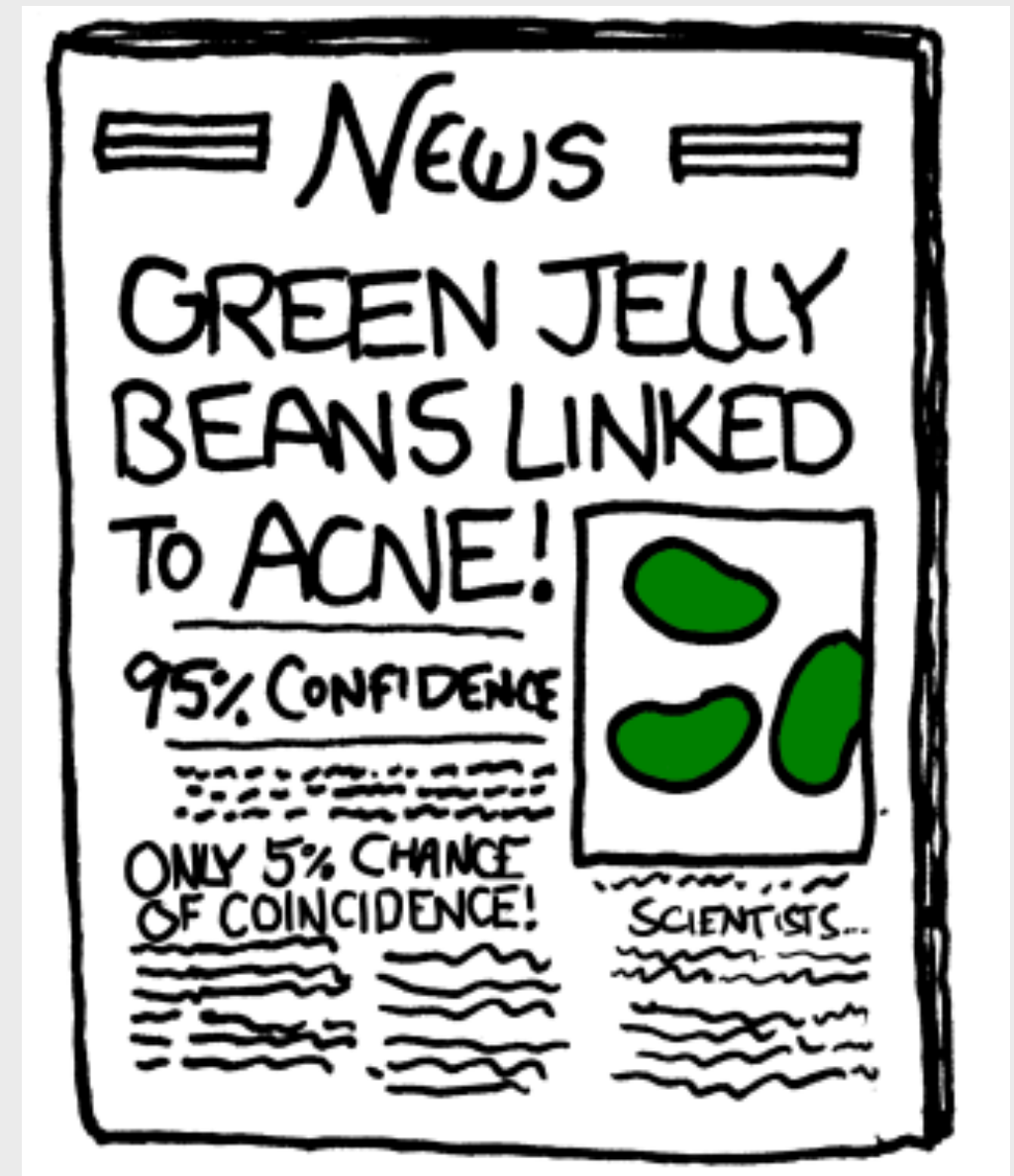
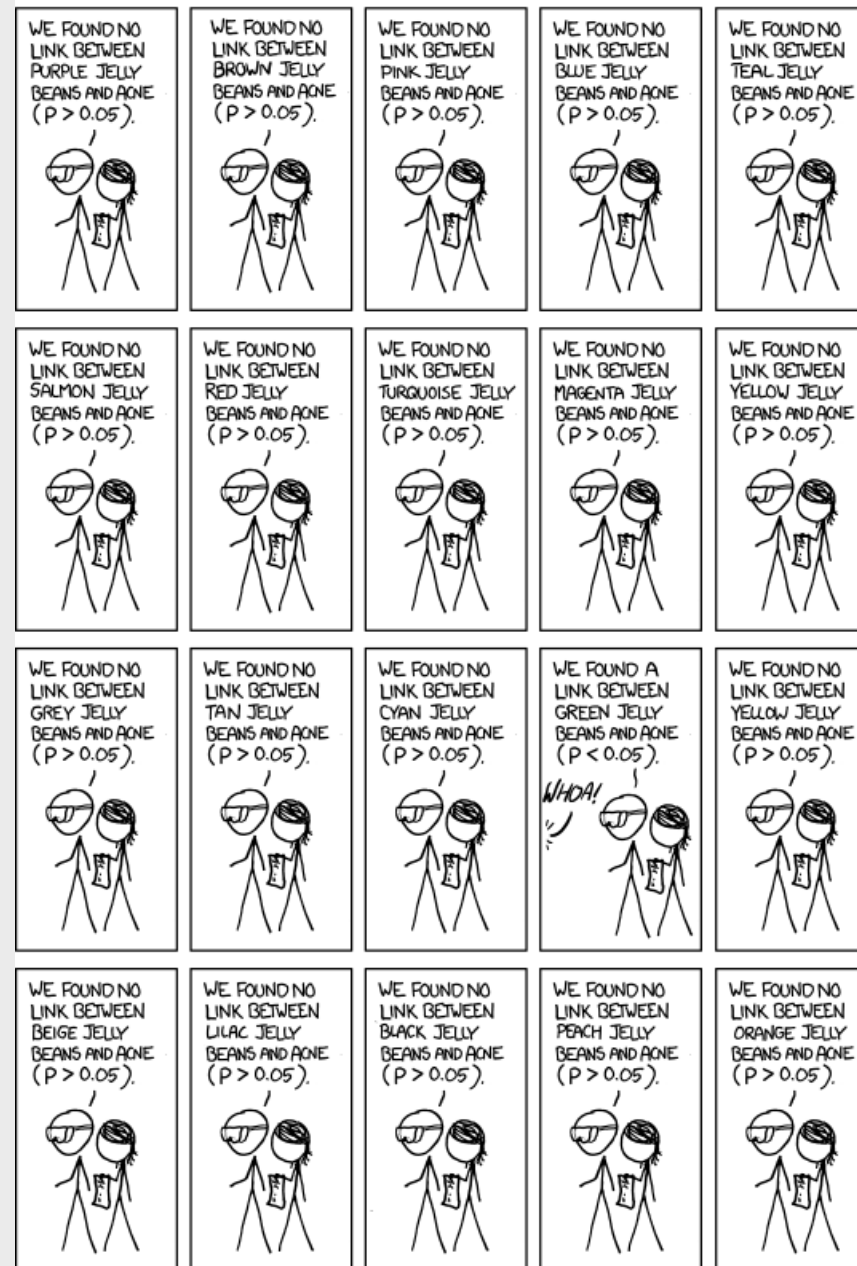
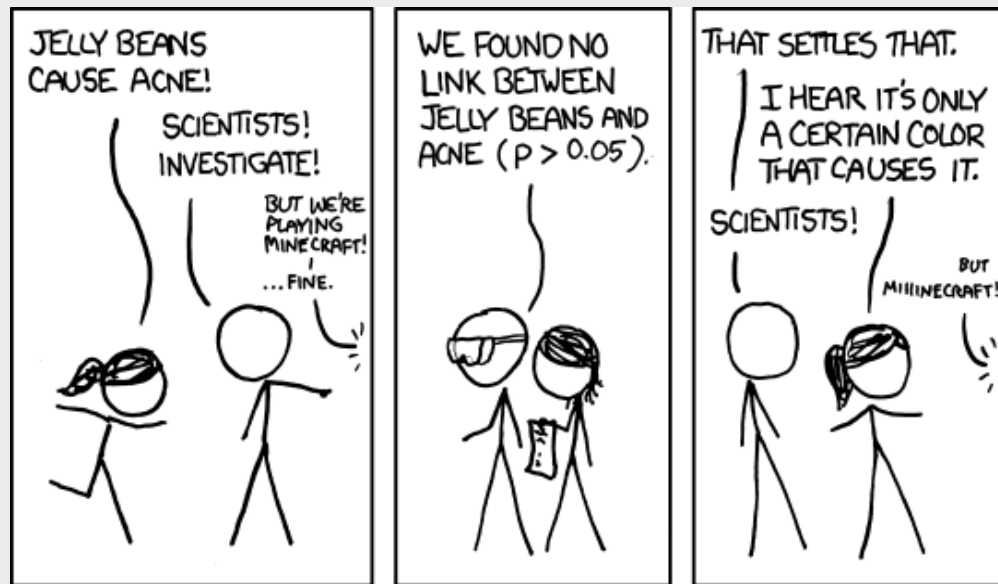
- **1-Sample:** \bar{x} vs. μ_0 . (eine Gruppe, Referenzwert).
- **2-Sample:** zwei unabhängige Gruppen.
- **Paired:** Vorher–Nachher innerhalb derselben Personen.
- Ziel: prüfen, ob Mittelwertsunterschiede Zufall sind.
- Immer p + CI + Effektgrösse berichten.

t-Tests zeigen, ob Mittelwertsunterschiede plausibel sind → aber erst Kontext und Effektgrösse machen sie bedeutsam.

Praxis & Kritik

p-Hacking: Wenn man das «p» jagt

- **Definition:** Wiederholtes Testen oder Subgruppen-Analysen, bis $p < 0.05$ erreicht ist.
- **Ziel:** ein „signifikantes“ Ergebnis finden, obwohl H_0 wahr ist.
- **Beispiel-Kaskade:** «Äpfel & Akne»
 - Test 1: alle Äpfel $\rightarrow p = 0.45$
 - Test 2: nur grüne $\rightarrow p = 0.08$
 - Test 3: nur Männer $\rightarrow p = 0.20$
 - ...
 - Test 20: gelbe Äpfel bei Teenagern $\rightarrow p = 0.04 \rightarrow$ scheinbar signifikant!
- **Problem:** Nur der letzte Test wird berichtet \rightarrow Falsch-Positiver Befund.



xkcd Comic: Significant

p-Hacking & Relevanz

- **p-Hacking:** wiederholtes Testen, bis $p < 0.05$.
- **Problem:** erhöht massiv die Wahrscheinlichkeit für Fehlalarme (α -Inflation).
- Immer dokumentieren, **wie viele Tests** durchgeführt wurden.

Signifikanz vs. Relevanz

- Statistisch signifikant \neq praktisch bedeutsam.
- Beispiel: Puls sinkt um 0.3 bpm, $p = 0.001 \rightarrow$ irrelevant im Alltag.
- Deshalb immer:
 - Effektgrösse (z. B. Cohen's d).
 - Konfidenzintervall.
 - Kontext und Bedeutung.

Signifikanz sagt, dass etwas auffällt \rightarrow nicht, dass es wichtig ist.

Best Practices

- Hypothesen **vorab registrieren** (Pre-Registration).
- Effektgrößen & **Konfidenzintervalle** immer berichten.
- Bei vielen Tests: **Bonferroni oder FDR** (Statcoach) anwenden.
- Kontext erläutern: Statistik liefert Evidenz, kein Urteil.
- Ziel: **Reproduzierbare Forschung**, nicht überraschende p-Werte.

Take-away: Praxis & Kritik

- **p-Hacking**: viele Tests → mehr Zufallstreffer.
- **Signifikanz \neq Relevanz**: immer Effektgrösse & Kontext prüfen.
- **Best Practice**: Pre-Registration, CI, Transparenz.

Klar planen, sauber berichten, ehrlich interpretieren.

Reflexion & Take-aways

Was wir heute gelernt haben

- **Zweck von Hypothesentests:** Zufall von echtem Effekt trennen.
- **Kernkonzepte:** H_0 , H_1 , α , β , Power, p-Wert.
- **t-Test:** prüft Mittelwertsunterschiede – Beispiel für angewandte Inferenz.
- **p-Wert:** misst Seltenheit unter H_0 , kein Wahrheitsmass.
- **Kritik & Praxis:** Signifikanz \neq Relevanz, p-Hacking vermeiden, Transparenz fördern.

Transferfrage: Interpretation üben

Eine Studie findet $p = 0.02$ und eine Effektgrösse $d = 0.1$.

Was bedeutet das?

- Statistisch **signifikant** (selten unter H_0).
 - Aber **praktisch irrelevant** (winziger Effekt).
- **Beispiel für „signifikant, aber nicht bedeutsam“.**

Quiz: *Aktive Wiederholung*

Kahoot Quiz VL7: Hypothesentests

Nächste Woche: Hypothesentests II

- Erweiterung auf komplexere Designs
 - [A/B-Tests](#)
 - [Permutationstests](#)
 - [ANOVA](#)
 - Fokus auf Effektgrößen, Power-Analyse und Multiple Testing.
- 👉 Vorbereitung: [PSDS Kapitel 4](#)