

M24 Statistik 1: Sommersemester 2024

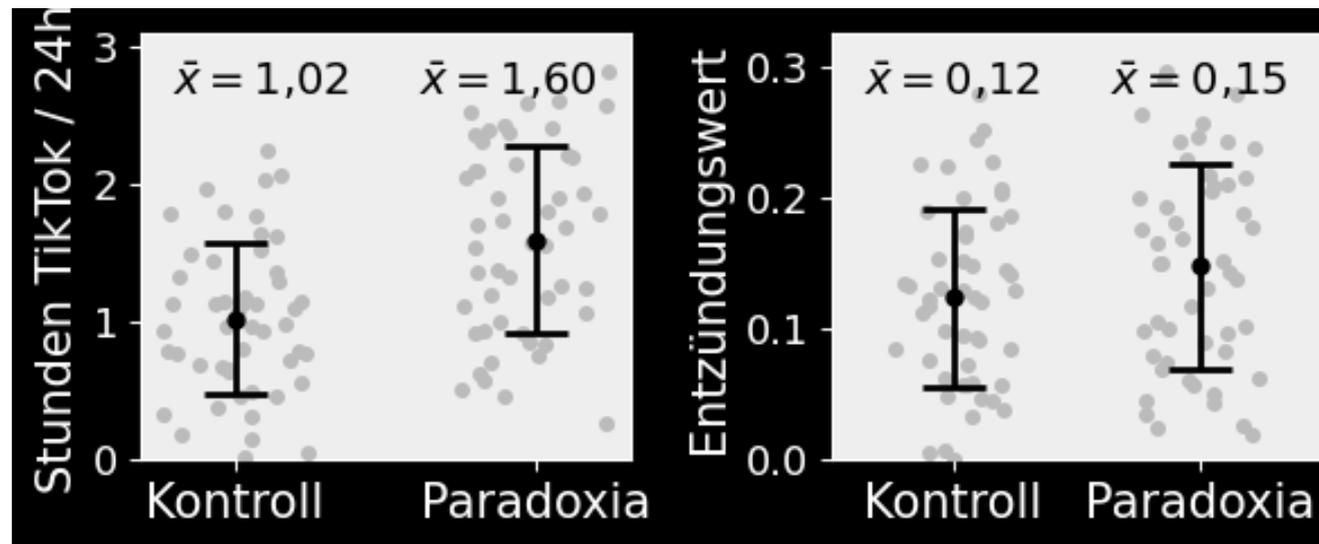
# Vorlesung 07: Effektstärke

Prof. Matthias Guggenmos

Health and Medical University Potsdam



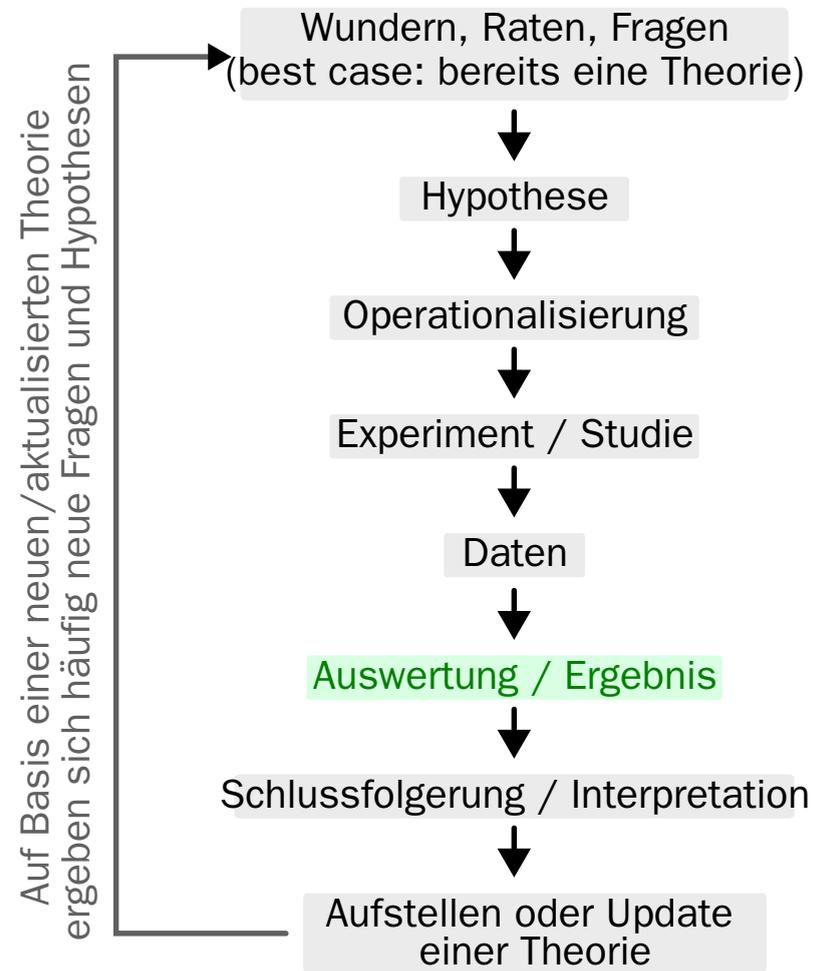
Sie sinnieren weiterhin über das Ergebnis Ihrer Beobachtungsstudie. Paradoxiker verbringen sowohl mehr Zeit auf TikTok, als auch weisen sie höhere Entzündungswerte auf. Zwar ist der Mittelwertsunterschied bei der TikTok-Zeit größer, aber Sie wissen, dass TikTok-Zeit und Entzündungswerte völlig unterschiedliche Skalen und daher nicht vergleichbar sind.



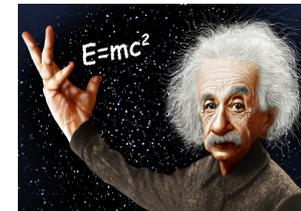
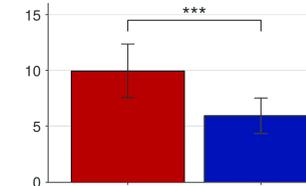
Die Frage lautet: wie kann man die beiden Gruppenunterschiede bezüglich TikTok-Zeit und Entzündungsparametern vergleichbar machen? Wie können wir eine Aussage treffen, welcher der beiden Effekte stärker ist?

# Effektstärke

# Der Forschungsprozess

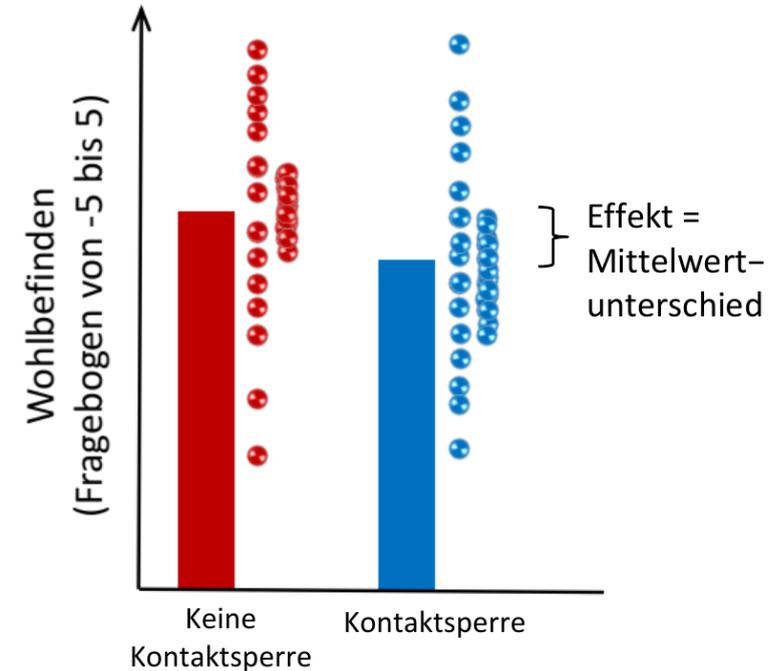


Player	Minutes	Points	Rebounds
A	41	20	6
B	30	29	7
C	22	7	7
D	26	3	3
E	20	19	8



# Effektstärke

- In psychologischer Forschung untersuchen wir in den meisten Fällen die Auswirkung von Variablen  $X_i$  auf Variablen  $Y_i$ .
- Diese Auswirkung ist entweder als **Unterschied** (z.B. wenn  $X$  die Gruppenzugehörigkeit angibt) oder als **Zusammenhang** (wenn  $X$  und  $Y$  metrische Variablen mit einer vermuteten kausalen Wechselwirkung sind) messbar – man spricht auch von **Effekten**.
- Wie kann die Aussagekraft bzw. Bedeutsamkeit von Effekten bestimmt und kommuniziert werden?
  - ⇒ **statistische Signifikanz** (ab Vorlesung 10): kann ein Effekt *allein durch Zufall erklärt werden*?
  - ⇒ **praktische Signifikanz** (Effektstärken): ist die Stärke des Effektes *in der Praxis bedeutsam*?
- Die Stärke eines Effektes im Sinne der praktischen Signifikanz wird als **Effektstärke** oder **Effektgröße** bezeichnet. Wir werden nachfolgend den Begriff Effektstärke verwenden.
- Unterschiedlichen Maße für die Effektstärke werden als **Effektmaße** bezeichnet.



Beispiel: Studie zum Wohlbefinden in Regionen mit und ohne Corona-Kontaktsperre

# Unstandardisierte und standardisierte Effektstärken

- Mittelwertsdifferenzen, Kovarianzen und Regressionskoeffizienten sind **unstandardisierte Effektstärken**, weil sie in den Rohwerten der Messung vorliegen.



**Beispiel Mittelwertunterschied:** der durchschnittliche Größenunterschied von erwachsenen Männern und Frauen in Deutschland beträgt 16cm<sup>1</sup>.



**Beispiel unstandardisierter Regressionskoeffizient:** je 0,1 Verbesserung in der **Abiturnote** steigt das monatliche Einstiegseinkommen um durchschnittlich 70 Euro<sup>2</sup>.

- In den beiden genannten Beispielen haben die Effektstärken sinnvolle und interpretierbare Einheiten und wären vergleichbar zwischen Studien.
- Gerade in der Psychologie ist dies aber nicht immer gegeben:
  - Fragebögen: Punktzahlen hängen willkürlich vom Kodierungsschema und der Zahl der Items ab
  - Ratingskalen: Ratingskalen unterscheiden sich häufig (Wohlbefinden auf einer Skala von 0 bis 100%, Wohlbefinden auf einer Skala von -5 bis 5, usw.)
- Falls die Skala (z.B. der verwendete Fragebogen) neu oder wenig bekannt ist, wie soll dann der Effekt interpretiert werden? Wann kann er als groß und wann als klein gelten?

# Standardisierte Effektstärken

- Um Effektstärken unabhängig von der verwendeten Skala zu vergleichen, werden Effektstärken **standardisiert**.
- Die Transformation der **Standardisierung** haben wir bereits bei Zufallsvariablen kennengelernt: *Teilen durch die Standardabweichung*.
  - Dieser Ansatz lässt sich analog auch auf Effekte beziehen
  - Ein Beispiel für einen Effekt ist die Mittelwertdifferenz — hier gilt:  
standardisierte Effektstärke = Mittelwertdifferenz / Standardabweichung
  - Besteht der Effekt aus einem Produkt zweier Variablen (wie bei der Kovarianz), muss der Effekt durch die Standardabweichung beider Variablen geteilt werden, um die Einheiten herauszukürzen.
- In der Folge sind standardisierte Effekte einheitslos, da die zugrundegelegten Standardabweichungen die gleiche Einheit wie die Effekte haben.
- Zu beachten ist, dass das Teilen durch die Standardabweichung der resultierenden Effektstärke eine spezifische Interpretation zuschreibt: *ein Effekt wird als “stärker” gewertet, wenn der unstandardisierte Effekt groß ist im Vergleich zur zugrundegelegten Standardabweichung*.
- Es gibt zwei wesentliche Funktionen / Einsatzzwecke von Effektstärken:
  - Einordnung der Stärke eines Effektes in einer *einzelnen Studie*
  - Vergleichbarmachung von Effekten zwischen *verschiedenen Studien*

# Mittelwertunterschiede

# Cohen's d

- Es werden drei Fälle von Mittelwertunterschieden unterschieden:
  - Fall 1: eine Stichprobe + Einzelmessung: Differenz zwischen dem Mittelwert einer Messung und einem Referenzwert (z.B. IQ=100)
  - Fall 2: zwei Stichproben + unabhängige Messungen: Differenz der Gruppenmittelwerte (z.B. IQ in einer blonden versus brünetten Gruppe)
  - Fall 3: eine Stichprobe + abhängige Messungen: Differenz der Mittelwerte zweier Messungen in derselben Gruppe (z.B. IQ morgens und IQ abends)
- Die standardisierte Effektstärke für alle drei Fälle eines Mittelwertunterschieds berechnet sich als

$$d = \frac{\text{Mittelwertdifferenz}}{\text{Standardabweichung}}$$

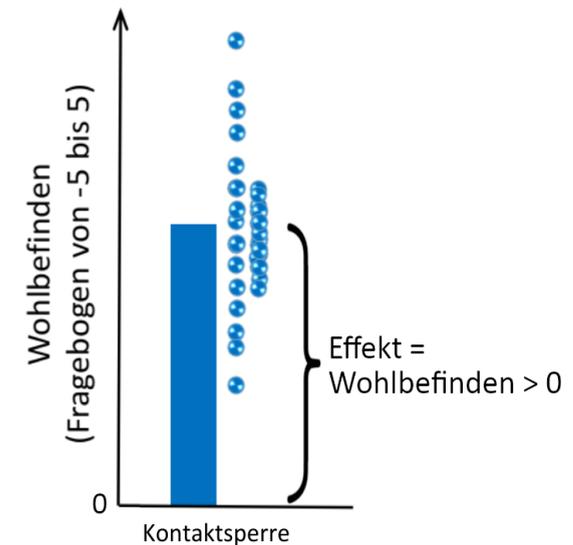
- Die Bezeichnung  $d$  für standardisierte Mittelwertunterschiede stammt von dem Statistiker Jacob Cohen — häufig wird daher auch von **Cohen's d** gesprochen.
- Während die Mittelwertdifferenz eindeutig ist, ist die weniger triviale Frage: Standardabweichung *von was?*

# Fall 1: eine Stichprobe + Einzelmessung

Es gibt nur eine einzelne Messung an einer Gruppe und die Frage ist, ob der Mittelwert  $\bar{x}$  bedeutsam über einem Referenzwert  $\mu_0$  liegt.

$$\text{Mittelwertdifferenz} = \bar{x} - \mu_0$$

Standardisierte Effektstärke: 
$$d = \frac{\bar{x} - \mu_0}{\hat{\sigma}}$$



Beispiel: Studie zum Wohlbefinden – ist das Wohlbefinden *in der Gruppe* mit Kontaktsperrre noch im positiven Bereich ( $\bar{x} > \mu_0$  mit  $\mu_0 = 0$ )?

- Die Wahl der Standardabweichung bereitet hier keine Kopfzerbrechen – es ist schlicht die Standardabweichung der Variable  $X$  in der Stichprobe.
- Eine kleine Besonderheit ist die Verwendung von  $n - 1$  statt  $n$  im Nenner der Varianzformel:

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

- Dies ist die sog. **Besselkorrektur**, die immer dann Anwendung findet, wenn auf Basis der Stichprobe eine Varianzschätzung für die Population getätigt werden soll – in aller Regel ist dies bei der Effektstärkenberechnung der Fall!

## Fall 2: zwei Stichproben + unabhängige Messungen

Es wird die Mittelwertdifferenz  $\bar{x}_A - \bar{x}_B$  zweier unabhängiger Gruppen berechnet und die Frage ist, ob diese Differenz bedeutsam ist.

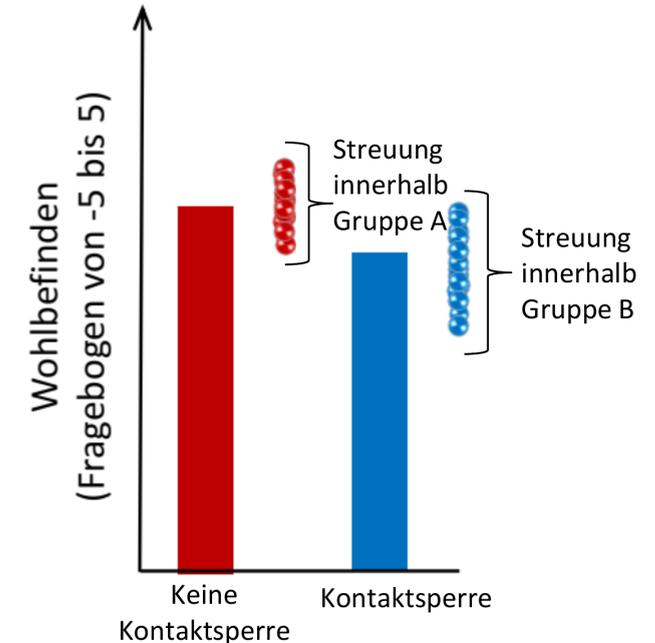
$$\text{Mittelwertdifferenz} = \bar{x}_A - \bar{x}_B$$

Standardisierte Effektstärke: 
$$d = \frac{\bar{x}_A - \bar{x}_B}{\hat{\sigma}_{\text{pooled}}}$$

- Im Nenner von Cohen's d wird hier die **gepoolte Varianz** verwendet.
- Es handelt sich dabei um die mittlere Varianz über beide Gruppen hinweg, entsprechend einem an den Stichprobengrößen  $n_A$  und  $n_B$  gewichteten Mittelwert:

$$\hat{\sigma}_{\text{pooled}} = \sqrt{\frac{(n_A - 1)\hat{\sigma}_A^2 + (n_B - 1)\hat{\sigma}_B^2}{n_A + n_B - 2}}$$

$$\text{Falls } n_A = n_B = n : \hat{\sigma}_{\text{pooled}} = \sqrt{\frac{\hat{\sigma}_A^2 + \hat{\sigma}_B^2}{2}}$$



Unterscheidet sich der Mittelwert zweier Gruppen?

## Einschub: gepoolte Varianz

- Bei unabhängigen Messungen ist es nicht (wie in Fall 3) möglich, die Einzelmesswerte  $x_{A,i}$  und  $x_{B,i}$  direkt voneinander abzuziehen. Nicht zuletzt wäre es völlig unklar, welche Versuchsperson aus A jeweils von welcher Versuchsperson aus B subtrahiert wird.
- Damit kann auch nicht die Varianz der individuellen Effekte der Versuchspersonen berechnet werden (“Differenzenvarianz”), sondern lediglich eine **gepoolte Varianz**  $\hat{\sigma}_{\text{pooled}}^2$ .
- Die Annahme ist dabei, dass beide Gruppen eine ähnliche Streuung haben, z.B. weil sie randomisiert aus derselben Population gezogen wurden und man nicht erwartet, dass die Gruppenzugehörigkeit die Varianz beeinflusst.
- Der Name leitet sich von der Vorstellung ab, die Datenpunkte beider Gruppen in einen gemeinsamen “Pool” zu werfen, und dann die Varianz auf allen Daten zu berechnen.
- Vor dem “Wurf in den Pool” müssen die Mittelwerte der Datenpunkte aus den jeweiligen Gruppen gleichgesetzt werden – im einfachsten Fall so, dass jeweils der Gruppenmittelwert von den Datenpunkten abgezogen wird (die Daten werden **zentriert** – beide Gruppen haben danach den Mittelwert 0):

$$X'_A = X_A - \bar{X}_A \quad \text{und} \quad X'_B = X_B - \bar{X}_B$$

$$X_{\text{pooled}} = [X'_A; X'_B]$$

$$\hat{\sigma}_{\text{pooled}}^2 = \hat{V}ar(X_{\text{pooled}})$$



# Einschub: gepoolte Varianz

- Es lässt sich mathematisch recht einfach zeigen, dass die gepoolte Varianz identisch dem an den Stichprobengrößen  $n_A$  und  $n_B$  gewichteten Varianzmittelwert ist:

$$\hat{\sigma}_{\text{pooled}}^2 = \frac{(n_A - 1)\hat{\sigma}_A^2 + (n_B - 1)\hat{\sigma}_B^2}{n_A + n_B - 2} \quad \text{bzw.} \quad \hat{\sigma}_{\text{pooled}} = \sqrt{\frac{(n_A - 1)\hat{\sigma}_A^2 + (n_B - 1)\hat{\sigma}_B^2}{n_A + n_B - 2}}$$

Die Subtraktion von 1 von den Stichprobengrößen  $n_A$  und  $n_B$  ist Ausdruck der Besselkorrektur.

- Die Gewichtung mit den Stichprobengrößen grenzt die gepoolte Varianz von der einfach gemittelten Varianz  $\hat{\sigma}_{\text{av}}^2$  ab:

$$\hat{\sigma}_{\text{av}}^2 = \frac{\hat{\sigma}_A^2 + \hat{\sigma}_B^2}{2} \quad \text{bzw.} \quad \hat{\sigma}_{\text{av}} = \sqrt{\frac{\hat{\sigma}_A^2 + \hat{\sigma}_B^2}{2}}$$

- $\hat{\sigma}_{\text{av}}^2$  wird auch als **ungepoolte Varianz** bezeichnet, weil hier nicht alle Datenpunkte in einen Pool geworfen werden, sondern die Varianzen  $\hat{\sigma}_A^2$  und  $\hat{\sigma}_B^2$  einzeln berechnet und anschließend ohne Berücksichtigung unterschiedlicher Stichprobengrößen gemittelt werden.



## Einschub: gepoolte Varianz

- Der Nachteil der gepoolten Varianz ist, dass sie nicht mehr der interindividuellen Variabilität des Effektes (hier Mittelwertdifferenz) entspricht, sondern der Variabilität des gemessenen Merkmals  $X$ .
  - Die Mittelwertdifferenz im Zähler wird also streng genommen nicht mehr ins Verhältnis zu ihrer eigenen Variabilität gesetzt.
- Die gepoolte Varianz hat jedoch auch einen Vorteil: es wird de facto nur eine Varianz geschätzt (die Populationsvarianz des Merkmals  $X$ ) und diese Schätzung basiert auf  $n_A + n_B$  Versuchspersonen. Sie kann also recht präzise geschätzt werden. Die Präzision der Varianz bestimmt wiederum direkt und wesentlich die Präzision des Effektstärkenmaßes selbst.
- Wenn Effektstärken in erster Linie zur Vergleichbarmachung mit anderen Studien dienen (Stichwort **Metaanalyse**), ist eine a) *präzise* und b) *zwischen Studien vergleichbare* Schätzung der Varianz ein wichtiger Aspekt — ggf. wichtiger als die Interpretierbarkeit des resultierenden Effektmaßes.



### Beispiel

Bei einem Pre-post-Interventionsdesign (Messung vor und nach einer Intervention an derselben Gruppe), kann es vorteilhaft sein, nur die *Standardabweichung des Vortestes* für die Standardisierung zu nehmen. Grund: die Variabilität *vor* einer Intervention ist am ehesten zwischen Studien vergleichbar — damit werden auch die Effektstärken besser vergleichbar.<sup>3</sup> In der Praxis geschieht dies aber selten (nicht immer ist die statistisch beste auch die populärste Praxis).





# Einschub: gepoolte Varianz

- Eine wichtige Voraussetzung für die Verwendung der gepoolten Varianz ist, dass die Einzelvarianzen  $\sigma_A^2$  und  $\sigma_B^2$  ähnlich sind.
- Warum? Die Idee ist, die Varianzschätzung des Merkmals durch Nutzung der Datenpunkte *beider* Gruppen zu verbessern. Sind die Varianzen in beiden Gruppen jedoch unterschiedlich, so ist der Fall gegeben, dass es gar nicht *die eine* Varianz des Merkmals gibt, sondern die Varianz stark von der Gruppenzugehörigkeit abhängt.
- Eine gängige Faustregel besagt, dass die Varianzen sich höchstens um den Faktor 2 unterscheiden sollten, d.h.  $0,5 < \frac{\sigma_A^2}{\sigma_B^2} < 2$ .
- Diese Einschränkung bei der Verwendung der gepoolten Varianz wird jedoch häufig missachtet, mutmaßlich auch, weil sich keine alternative Lösung für den Fall unabhängiger Messungen in großer Breite etabliert hat.
  - In einem aktuellen Preprint schlagen Delacre und Kollegen für den Fall unabhängiger Messungen und unterschiedlicher Varianzen die Effektstärkengröße *Hedges' g\** auf Basis der ungepoolten Varianz  $\hat{\sigma}_{av}$  vor<sup>4</sup>:

$$\text{Hedges' } g^* \approx \left( 1 - \frac{3}{4(n_A + n_B) - 9} \right) \cdot \frac{\bar{x}_A - \bar{x}_B}{\hat{\sigma}_{av}} \quad \text{mit} \quad \hat{\sigma}_{av} = \sqrt{\frac{\hat{\sigma}_A^2 + \hat{\sigma}_B^2}{2}}$$

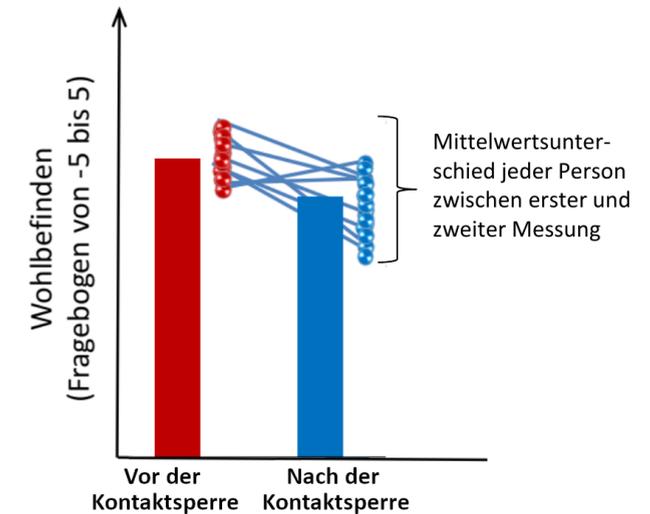
## Fall 3: eine Stichprobe + zwei abhängige Messungen

Es wird die Mittelwertdifferenz  $\bar{x}_A - \bar{x}_B$  zweier abhängiger Messungen in einer Stichprobe berechnet und die Frage ist, ob diese Differenz bedeutsam ist.

$$\text{Mittelwertdifferenz} = \bar{x}_A - \bar{x}_B$$

Standardisierte Effektstärke:

$$d = \frac{\bar{x}_A - \bar{x}_B}{\hat{\sigma}_\Delta} \quad \text{oder} \quad d = \frac{\bar{x}_A - \bar{x}_B}{\hat{\sigma}_{\text{pooled}}}$$



Beispiel: Unterscheidet sich das Wohlbefinden in derselben Gruppe vor und nach einer Kontaktsperre?

### Variante a): Differenzvarianz

- In einer gängigen Variante wird die Standardabweichung  $\hat{\sigma}_\Delta$  der Differenzvariable  $\Delta X = X_A - X_B$  gebildet:

$$\hat{\sigma}_\Delta = \sqrt{\frac{\sum (\Delta x_i - \Delta \bar{x})^2}{n - 1}} \quad \text{und} \quad \Delta x_i = x_i^{(A)} - x_i^{(B)}$$

- Wir bezeichnen diese als **Differenzvarianz**.



# Einschub: Differenzenvarianz

- Im Falle abhängiger Messungen können wir nicht nur die Mittelwerte  $\bar{x}_A$  und  $\bar{x}_B$  subtrahieren ( $\Delta\bar{x} = \bar{x}_A - \bar{x}_B$ ), sondern bereits die einzelnen Messwerte  $x_{A,i}$  und  $x_{B,i}$ :  $\Delta x_i = x_{A,i} - x_{B,i}$
- Die Differenzenvarianz misst die Varianz dieser Differenzwerte. Wir können sie mit der gewohnten Varianzformel berechnen, nur dass statt  $x_i$  die  $\Delta x_i$  eingesetzt werden:

$$\hat{\sigma}_{\Delta}^2 = \frac{1}{n-1} \sum (\Delta x_i - \Delta\bar{x})^2 \quad \text{bzw.} \quad \hat{\sigma}_{\Delta} = \sqrt{\frac{1}{n-1} \sum (\Delta x_i - \Delta\bar{x})^2}$$

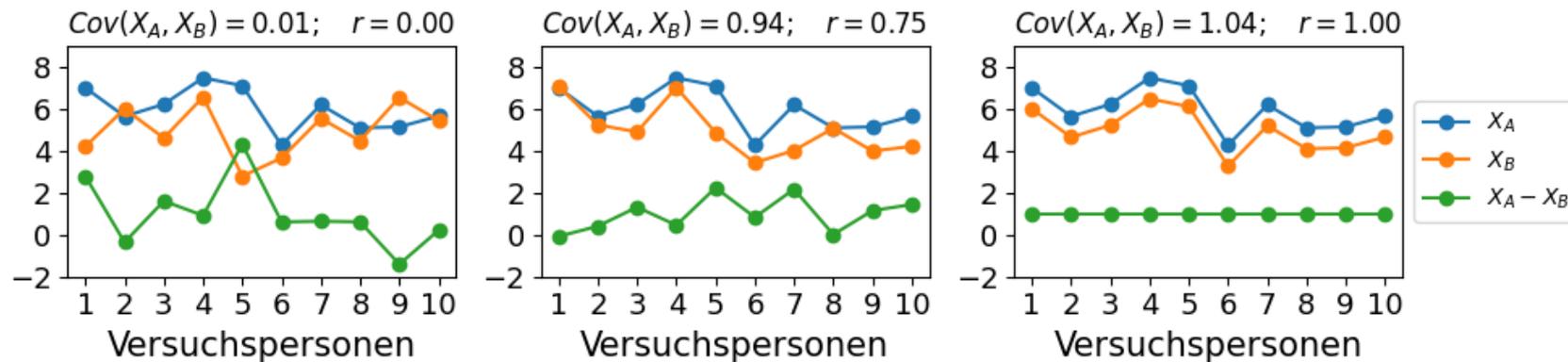
- Mit einigen mathematischen Tricks lässt sich zeigen, dass die Formel für die Varianz der Differenzwerte auch wie folgt dargestellt werden kann:

$$\hat{\sigma}_{\Delta}^2 = \hat{\sigma}_A^2 + \hat{\sigma}_B^2 - 2 \hat{Cov}(X_A, X_B)$$

- Diese Formel macht transparent, dass die Differenzenvarianz vom Zusammenhang zwischen  $X_A$  und  $X_B$  abhängt:
  - Sind die Zufallsvariablen *nicht korreliert* ( $\hat{Cov}(X_A, X_B) = 0$ ), so ist die Varianz der Differenz einfach der Summe der Einzelvarianzen.
  - Sind die Zufallsvariablen *positiv korreliert* ( $\hat{Cov}(X_A, X_B) > 0$ ), so reduziert sich die Summe in Abhängigkeit von der Kovarianz.
  - Sind die Zufallsvariablen *negativ korreliert* ( $\hat{Cov}(X_A, X_B) < 0$ ), so erhöht sich die Summe in Abhängigkeit von der Kovarianz.

# Einschub: Differenzenvarianz

Der Effekt der Kovarianz auf die Varianz der Differenzwerte ist besonders intuitiv bei einer positiven Korrelation. Betrachten wir zwei Variablen  $X_A$  und  $X_B$  mit unterschiedlichen Kovarianzen und wie sich dabei jeweils die Variabilität des Differenzwertes entwickelt:



In allen drei Plots gilt  $\bar{x}_A = 6$  und  $\bar{x}_B = 5$ , d.h.  $\bar{x}_A - \bar{x}_B = 1$ .

Wir sehen: je höher die Korrelation  $\bar{x}_A$  und  $\bar{x}_B$ , desto geringer die Variabilität der Differenz von  $\bar{x}_A - \bar{x}_B$ ! Ist die Korrelation perfekt ( $r = 1$ ), so ist die Differenz sogar konstant, d.h. ihre Variabilität ist gleich 0.



## Fall 3: eine Stichprobe + zwei abhängige Messungen

- Bei der Standardisierung mit  $\hat{\sigma}_{\Delta}^2$  wird der Effekt (Mittelwertdifferenz) mit der *Streuung des Effektes* (Varianz der Mittelwertdifferenz) ins Verhältnis gesetzt.
- Diese Art der Standardisierung entspricht der intuitiven Grundidee von Cohen's d: den unstandardisierte Effekt an seiner eigenen Streuung zu relativieren (d.h. zu standardisieren).

### Variante b): gepoolte Varianz

- **Problem:** häufig sollen Effektstärken zwischen Designs mit *abhängigen* und *unabhängigen* Messungen verglichen werden — jedoch lässt sich in Fall 2 (unabhängige Messungen) keine Differenzenvarianz berechnen.
- Aus diesem Grund gibt es eine zweite gängige von Cohen's d bei abhängigen Messungen, die auf der **gepoolten Varianz** basiert (analog zu Fall 2):

$$\hat{\sigma}_{\text{pooled}}^2 = \frac{\hat{\sigma}_A^2 + \hat{\sigma}_B^2}{2} \quad \text{bzw.} \quad \hat{\sigma}_{\text{pooled}} = \sqrt{\frac{\hat{\sigma}_A^2 + \hat{\sigma}_B^2}{2}}$$

- Beachte, dass hier die gepoolte Varianz  $\hat{\sigma}_{\text{pooled}}^2$  identisch mit der ungepoolten (einfach gemittelten) Varianz  $\hat{\sigma}_{\text{av}}^2$  ist, da bei abhängigen Messungen immer gilt:  $n = n_A = n_B$ .

# Aktuelle Forschung

- Laut einer aktuellen Forschungsarbeit<sup>5</sup> sind die Werte von Cohen's d auf Basis der gepoolten Varianz nicht exakt vergleichbar zwischen abhängigen und unabhängigen Messungen (auch dann, wenn die Varianzen ähnlich sind).
- Die Arbeit bietet dafür folgende Modifikation von Cohen's d für abhängige Messungen an:

$$d = \sqrt{\frac{2(1-r)}{n}} \cdot t'_{\nu}(\lambda) \quad \text{mit} \quad \lambda = \sqrt{\frac{n}{2(1-r)}} \cdot \frac{\bar{x}_A - \bar{x}_B}{\hat{\sigma}_{\text{pooled}}}$$

wobei  $r$  die Pearson-Korrelation zwischen den beiden abhängigen Messungen ist und  $t'$  die nichtzentrale  $t$ -Verteilung mit  $\nu = 2(n-1)/(1+r^2)$  Freiheitsgraden. Beachte, dass die Formel nur gilt, falls ähnliche Varianzen in den beiden Bedingungen A und B angenommen werden können.

- Laut dem Autor Denis Cousineau der Forschungsarbeit ist damit Cohen's d exakt vergleichbar zwischen abhängigen und unabhängigen Messungen.
- Weiterer nützlicher Link zu Effektstärken bei abhängigen Messungen: <sup>6</sup>



# Interpretation

# Interpretation von Cohen's $d$

- Die Werte von Cohen's  $d$  reichen von  $-\infty$  bis  $+\infty$ . Negative Werte sind möglich, da das Vorzeichen davon abhängt, welcher Mittelwert von welchem abgezogen wird.
  - Es ist Konvention,  $d$  so zu berechnen, dass  $d$  positiv ausfällt, falls der Effekt in die hypothesierte Richtung geht.
  - Die Interpretation der Effektstärke macht sich aber am absoluten Wert  $|d|$  fest: wenn  $d = -0,3$ , dann hat der Effekt eine Stärke von  $d = 0,3$ , aber in eine andere Richtung als erwartet.
- Durch die Standardisierung mit der Standardabweichung gilt umgekehrt, dass Cohen's  $d$  ausdrückt, um wie viel Standardabweichungen ein Effekt von einem Nulleffekt abweicht. Diese Interpretation ist am intuitivsten für den Fall einer *Einzelmessung mit Referenzwert*:



Beispiel

Eine Studie untersucht, ob die Schlafdauer in einer Gruppe von Psychologiestudierenden in der Bachelorarbeitsphase geringer ist als die durchschnittliche Schlafdauer in Deutschland (7:45 Stunden). Tatsächlich zeigt sich eine verringerte Schlafdauer mit einem Cohen's  $d$  von 0,3.

- Die Effektstärke von 0,3 sagt aus, dass die Schlafdauer um 0,3 *Standardabweichungen* gegenüber dem Durchschnittswert verringert ist. Die Einheit *Standardabweichung* bezieht sich dabei auf die Standardabweichung  $\hat{\sigma}$  der Schlafdauer in der untersuchten Gruppe.
  - Bei mehreren Bedingungen/Gruppen gilt zwar weiterhin die Interpretation im Sinne von *in Einheiten von Standardabweichungen*, allerdings ist die Definition der Standardabweichungen ( $\hat{\sigma}_{\Delta}$ ,  $\hat{\sigma}_{\text{pooled}}$ ) eher kompliziert und damit nicht mehr sonderlich intuitiv.

# Interpretation von Effektstärken

- Um Effektstärken besser einordnen und kommunizieren zu können, hat Jacob Cohen folgende Unterteilung vorgeschlagen:

d	r	Interpretation
< 0.2	< 0.1	Trivialer Effekt
ab 0.2	ab 0.1	Kleiner Effekt
ab 0.5	ab 0.3	Mittlerer Effekt
ab 0.8	ab 0.5	Großer Effekt

- Zugleich fügt aber Cohen selbst folgende Qualifizierung an:

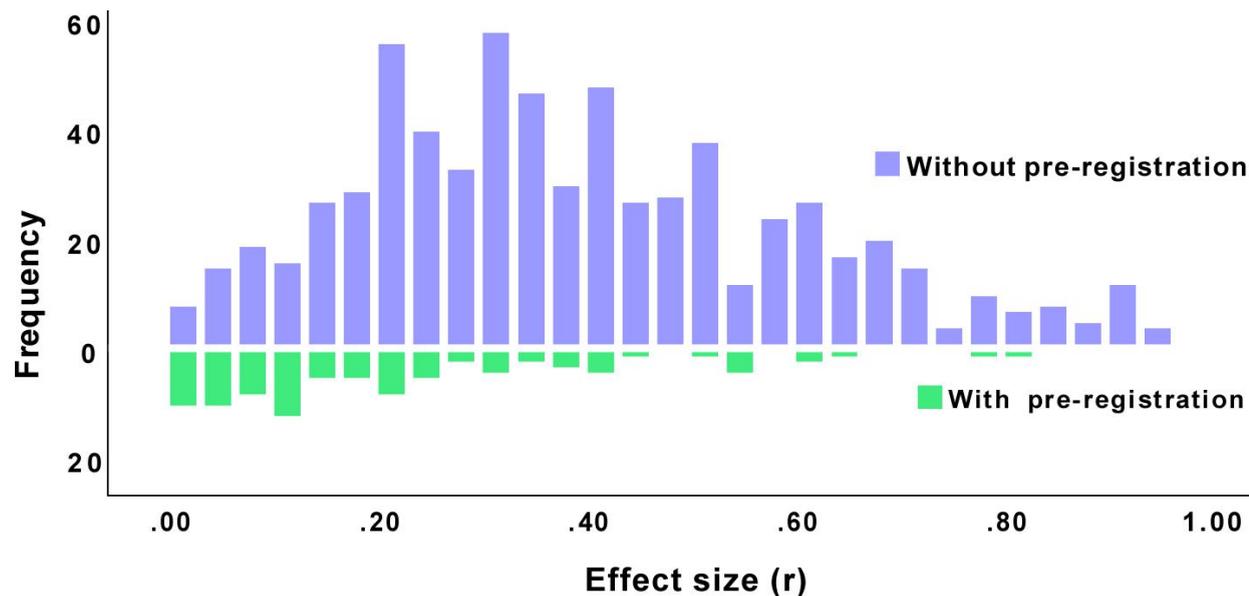


The terms „small," "medium," and „large" are relative, not only to each other, but to the area of behavioral science or even more particularly to the specific content and research method.  
(Cohen, 1988, p. 25)

- Effektstärken sollten also idealerweise in ihrem jeweiligen Kontext interpretiert werden.
- Beispiel:** Effektstärken bezüglich der Veränderung des Körpergewichts durch Diäten sind erwartbar größer, als Veränderungen bei eher stabilen Merkmalen wie Persönlichkeitsfacetten.  $\Rightarrow$  Ein d-Wert der einer vergleichsweise geringen Veränderung des Körpergewichts entspricht, würde vielleicht in der Persönlichkeitsforschung als starker Effekt gelten.

# Interpretation von Cohen's d

- Um zu beurteilen, was als kleine / mittlere / starke Effekte in einem spezifischen Kontext gilt, konsultiert man prinzipiell die entsprechende Fachliteratur nach typischen Referenzeffektstärken.
- Problem: es gibt gute Evidenz, dass *publizierte Effekte die wahren Effekte überschätzen (Publikationsbias)* ⇒ führt zu falschen Maßstäben



Die Abbildung zeigt, dass Effekte, bei denen Hypothesen und Analysen vorab registriert wurden (“with pre-registration”) deutlich geringere Effektstärken aufweisen, als Effekte “without pre-registration”. In der Abbildung ist zu beachten, dass die absolute und nicht die relative Häufigkeit aufgetragen ist, wodurch Studien ohne Preregistrierung – von denen es deutlich mehr gibt – visuell dominieren.<sup>7</sup>

⇒ Die Interpretation und Einordnung von Effektstärken ist ein nicht-triviales Problem, das viel “domain knowledge” erfordert. Dies gilt für standardisierte wie unstandardisierte Effektstärken.

- Faustregeln finden sich hier: <https://imaging.mrc-cbu.cam.ac.uk/statswiki/FAQ/effectSize>

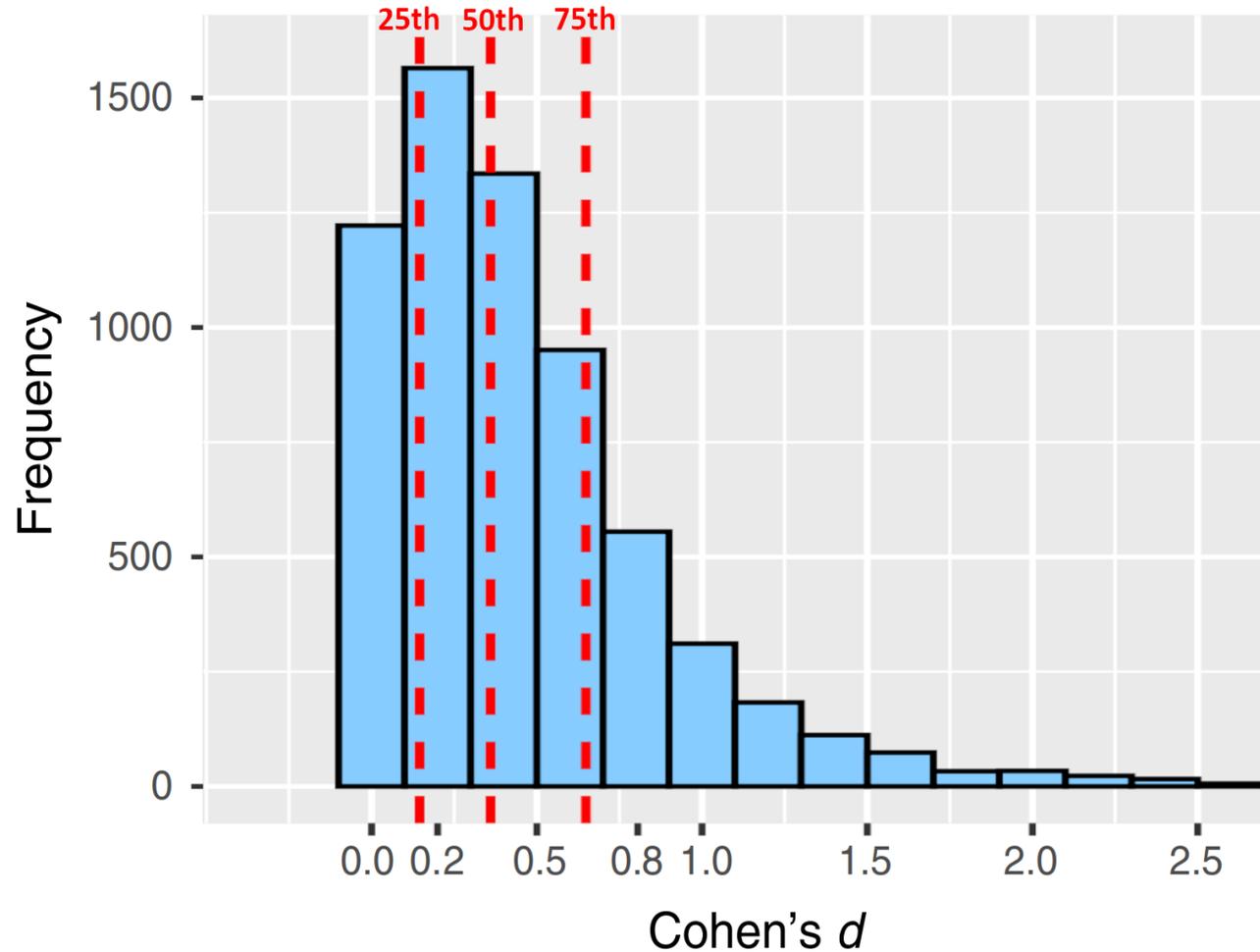
# Typische Effektstärken in der Sozialpsychologie

Subgroup	Number of meta-analysis	Number of effect sizes	Median	Mean	SD
<i>Correlation</i>					
Groups	15	998	0.26	0.31	0.25
Interpersonal relationships	12	2,323	0.30	0.32	0.19
Prejudice	10	2,639	0.18	0.21	0.15
Self	10	1991	0.29	0.31	0.20
Attitude	9	2,352	0.26	0.29	0.20
Social cognition	9	1,248	0.27	0.33	0.24
Gender differences	5	585	0.23	0.27	0.20
<i>Cohen's d</i>					
Gender differences	12	1,261	0.22	0.30	0.31
Prejudice	10	1,370	0.34	0.44	0.40
Self	10	884	0.48	0.59	0.56
Interpersonal relationships	9	1,075	0.28	0.39	0.41
Social cognition	9	750	0.50	0.58	0.52
Attitude	5	428	0.39	0.47	0.38

Typische Effektstärke in der sozialpsychologischen Literatur<sup>8</sup>.



# Typische Effektstärken in der Sozialpsychologie



Trotz Publikationsbias sind die tatsächlich in der Literatur berichteten Effektstärken geringer als bei der Einteilung nach Cohen angenommen. Auf Basis dieser Studie wären die korrekten Grenzen für Cohen's  $d$  0,15, 0,36, und 0,65. Bild adaptiert von Lovakov & Agadulina (2021)<sup>9</sup>.





# Interpreting Cohen's $d$ Effect Size

## An Interactive Visualization

Created by [Kristoffer Magnusson](#)

Share



The Cohen's  $d$  effect size is immensely popular in psychology. However, its interpretation is not straightforward and researchers often use general guidelines, such as small (0.2), medium (0.5) and large (0.8) when interpreting an effect. Moreover, in many cases it is questionable whether the standardized mean difference is more interpretable than the unstandardized mean difference.

In order to aid the interpretation of Cohen's  $d$ , this visualization offers these different representations of Cohen's  $d$ : visual overlap, Cohen's  $U_3$ , the probability of superiority, percentage of overlap, and the number needed to treat. It also lets you change the standard deviation and displays the

# Standardisierte oder unstandardisierte Effektstärken?

- Die Frage ob Effektstärken in standardisierter oder unstandardisierter Form berichtet werden sollten wird durchaus kontrovers diskutiert<sup>10 11</sup>.
- Standardisierte Effektstärken ermöglichen eine bessere Vergleichbarkeit zwischen unterschiedlichen Skalen, gleichzeitig wird die intuitive Bedeutung von *Effektstärke* aber verwässert.



Beispiel

Es wird berichtet, dass ein Coronaimpfstoff die Viruslast bei einer Infektion reduziert. Die Effektstärke wird mit Cohen's  $d = 0.4$  angegeben.

- Aus diesem Beispiel wird klar, dass die Effektstärke zum einen wenig intuitiv ist (in jedem Fall für Nicht-Wissenschaftler:innen), zum anderen ist nicht ersichtlich, ob die Viruslast nennenswert reduziert wurde oder ob der Rückgang eher klein war, aber die Standardabweichung in der Gruppe so gering, dass dennoch ein hoher  $d$ -Wert erreicht wurde.
- Darüber hinaus hängt die Standardabweichung einer Variable häufig mit eher nebensächlichen Details eines experimentellen Designs (within-subject vs. between-subject) oder einer Stichprobe (nur Psychologiestudierende oder heterogeneres Sample der Allgemeinbevölkerung?) zusammenhängt, die für die Effektstärke wenig relevant sind.
- Aus diesen Gründen sollte für **Effekte, die in interpretierbaren Einheiten vorliegen, immer (auch) die unstandardisierte Effektstärke** angegeben werden (z.B. Notenstufen, Einkommen, IQ-Punkte, Größe- Gewichtsangaben, Zeitangaben).

# Effektstärke bei Zusammenhängen

# Korrelationskoeffizient

- Im Fall von Zusammenhangsanalysen haben wir die standardisierte Effektstärke bereits kennengelernt: der **Korrelationskoeffizient** ( $r, r_s, \tau, \phi$ ). Beispiel Pearson-Korrelation:

$$r = \frac{Cov(X, Y)}{s_X s_Y}$$

- Der Nenner  $s_X s_Y$  stellt hier die Standardisierung dar.
- Wir können sehen, dass  $r$  standardisiert ist, da es keine Einheit hat und eine vergleichbar ist zwischen verschiedenen Skalen und verschiedenen Variablen.



Alle Korrelationskoeffizienten ( $r, r_s, \tau, \phi$ ) sind bereits Effektstärken.

---

# Umrechnung von Cohen's $d$ und Korrelationskoeffizient $r$

Fassen Metaanalysen sowohl Studien zusammen, die Effekte als Korrelation berichten (Effektmaß  $r$ ), als auch Studien, die Effekte als Mittelwertsunterschiede zwischen Bedingungen/Gruppen berichten (Effektmaß  $d$ ), entsteht die Notwendigkeit  $r$  und  $d$  ineinander umzurechnen.

**Fall 1: eine der Variablen in der Korrelation ist eine natürliche binäre Variable (z.B. männl./weibl.)**

In diesem Fall gilt:

$$d = \frac{2r}{\sqrt{1 - r^2}} \quad \hat{se}(d) = \frac{2}{\sqrt{(n - 1)(1 - r^2)}}$$

wobei  $\hat{se}(d)$  der Standardfehler der Effektstärke  $d$  ist, der zusätzlich angegeben werden sollte.

Umgekehrt gilt:

$$r = \frac{d}{\sqrt{d^2 + 4}}$$



# Umrechnung von Cohen's $d$ und Korrelationskoeffizient $r$

Fall 2: beide Variablen in der Korrelation sind kontinuierlich, oder eine Variable ist binär, aber entstand durch Dichotomisierung einer kontinuierlichen Variable (advanced!)

In diesem Fall muss eine Variable als unabhängige Variable  $X$  definiert werden. Falls eine der Variablen durch Dichotomisierung binär ist, ist diese Variable in jedem Fall die unabhängige Variable.

Es gilt:

$$d = \frac{kr}{\hat{\sigma}_X \sqrt{1 - r^2}} \quad \hat{se}(d) = |d| \sqrt{\frac{1}{r^2(n-3)} + \frac{1}{2(n-1)}}$$

wobei  $\hat{se}(d)$  der Standardfehler der Effektstärke  $d$  ist, der zusätzlich angegeben werden sollte.

**Interpretation:**  $d$  entspricht der durchschnittlichen Zunahme der standardisierten  $Y$ -Variable mit jeder Zunahme von  $X$  um  $k$  (Rohwert)Einheiten –  $k$  muss vom Forscher gewählt werden.

Hier wird klar, dass die Formel in Fall 1 implizit  $k = 2\hat{\sigma}_X$  annimmt. Wählt man dieses  $k$  auch in Fall 2, ist die Berechnung von  $d$  identisch zu Fall 1 – der Standardfehler unterscheidet sich allerdings weiterhin! Quelle: <sup>12</sup>

Umgekehrt gilt:

$$r = \frac{d\hat{\sigma}_X}{\sqrt{d^2\hat{\sigma}_X^2 + k^2}}$$



# Weitere Effektmaße

# Effektstärken für Mittelwertunterschiede bei mehr als zwei Messungen



- Gibt es mehr als zwei Experimentalbedingungen oder Gruppen (A, B, C, ...), gibt es zwei Möglichkeiten:
  1. Von Interesse sind die **paarweisen** Mittelwertunterschiede (z.B:  $A - B$ ,  $A - C$ ,  $B - C$ ) → in diesem Fall kann wie bisher Cohen's d für jedes Paar angewendet werden.
  2. Von Interesse ist, ob sich die Mittelwerte in den Gruppen A, B, C **in ihrer Gesamtheit betrachtet** unterscheiden, d.h. ob die Aufteilung in diese spezifischen Gruppen sinnvoll ist.

Fall 2 ist unser erster Kontakt mit der **Varianzanalyse**, die in Statistik 2 ausführlich behandelt wird. Man kann die Fragestellung in Fall 2 auch folgendermaßen formulieren:

Zu welchem Grad wird die Varianz der gepoolten Daten aller Gruppen bereits erklärt durch die Mittelwerte der Gruppen?

- Auf Basis dieser Formulierung ist nicht mehr überraschend, dass die Effektstärke für Mittelwertunterschiede von mehreren Messungen als *Verhältnis zweier Streuungen* ausgedrückt werden kann:

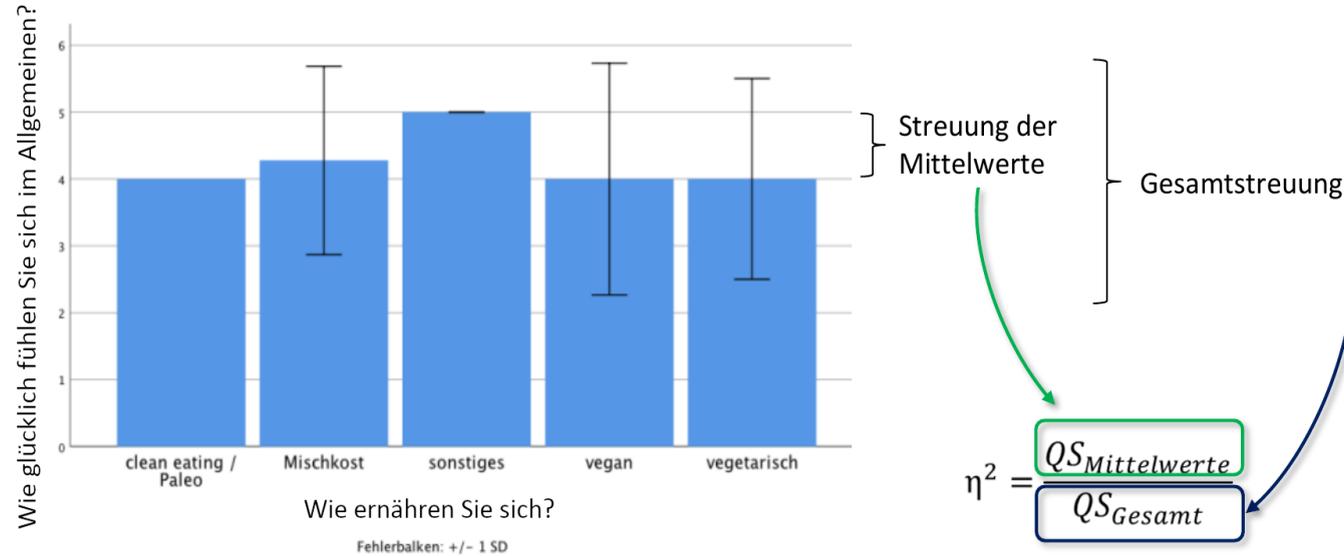
$$\eta^2 = \frac{QS_{\text{Mittelwerte}}}{QS_{\text{Gesamt}}}$$

$$\sigma^2 = \frac{\sum (x_i - \bar{x})^2}{n}$$

Quadratsumme (QS)

Die Quadratsumme entspricht dem Zähler in der Formel für die Varianz.

# Effektstärken für Mittelwertunterschiede bei mehr als zwei Messungen



Courtesy of Prof. Thomas Schäfer, Medical School Berlin

- $\eta^2$  ("Eta Quadrat") gibt an, wie viel der Gesamtvarianz durch die Varianz der Mittelwerte aufgeklärt wird.
- Es kann zwischen 0 (Mittelwerte erklären keine Varianz) und 1 (Mittelwerte erklären die komplette Varianz) liegen.
- Die Berechnung von  $QS_{\text{Gesamt}}$  umfasst 1) die Varianz zwischen den Bedingungen, 2) die Varianz zwischen Personen (über Bedingungen hinweg), und 3) wie sehr sich die Varianz der Bedingungen zwischen Personen unterscheidet:

$$QS_{\text{Gesamt}} = QS_{\text{Bedingungen}} + QS_{\text{Personen}} + QS_{\text{Personen} \times \text{Bedingungen}}$$



# Effektstärken für Mittelwertunterschiede bei mehr als zwei Messungen

- Es kann argumentiert werden, dass die Varianz, die lediglich die interindividuellen Unterschiede der Versuchspersonen (Varianzanteil 2 bzw.  $QS_{\text{Personen}}$ ) charakterisiert, für die Effektstärke irrelevant ist und nicht zu  $QS_{\text{Gesamt}}$  gezählt werden sollte, d.h.:

$$QS_{\text{Gesamt}} = QS_{\text{Bedingungen}} + QS_{\text{Personen} \times \text{Bedingungen}}$$

- Wird die Varianz  $QS_{\text{Personen}}$  nicht berücksichtigt, spricht man vom **partiellen  $\eta_p^2$** .
- Eine ausführliche Diskussion zum Pro und Kontra von  $\eta^2$  vs.  $\eta_p^2$  findet sich z.B. im Buch von Eid, Gollwitzer und Schmitt im Kapitel zur Varianzanalyse.



Nota Bene

Warum werden bei der Berechnung von  $\eta^2$  die Quadratsummen und nicht die Varianzen direkt verwendet? Grund ist, dass die Varianz die *durchschnittliche* (quadrierte) Abweichung vom Mittelwert angibt (daher der Faktor  $\frac{1}{n}$ ), und beim Vergleich verschiedener Variabilitätskomponenten nicht festgestellt werden kann, wie viel der Datenvariabilität *absolut gesehen* durch eine Variabilitätskomponente erklärt wird.



# Absolute Risikoreduktion (ARR)

Sind beide Variablen dichotom, sind weder der Korrelationskoeffizient noch Cohen's d intuitive Effektmaße.



Sie untersuchen, ob ein neues Medikament die Heilungsrate (Erfolgsrate) einer Krankheit verbessert. Die Treatmentgruppe erhält das Medikament, die Kontrollgruppe Placebo.

Eine sinnvolles Effektmaß ist hier, *um wie viel* die Erfolgsrate in der Treatmentgruppe die Erfolgsrate in der Kontrollgruppe übersteigt. Dies lässt sich einfach aus einer Vierfeldertafel ableiten:

	male	female	$\Sigma$
nonsmoker	a	b	a+b
smoker	c	d	c+d
$\Sigma$	a+c	b+d	n

$$ARR = \frac{a}{a+b} - \frac{c}{c+d}$$

Erfolgsrate in der Treatmentgruppe

Erfolgsrate in der Kontrollgruppe

- $ARR$  ist die **Absolute Risikoreduktion**.

# Numbers Needed to Treat (NNT)

- Noch gängiger als die Absolute Risikoreduktion ist die inverse Größe, die als **Numbers Needed to Treat (NNT)** bezeichnet wird.

---

**Definition** **Number Needed to Treat:** Anzahl der Personen, die behandelt werden müssten, damit eine zusätzliche Person einen Nutzen hat.

---

Mathematisch ist  $NNT$  das Inverse der  $ARR$ :

$$NNT = \frac{1}{ARR}$$

- Da es um Personen geht, wird die NNT immer aufgerundet.

 Beispiel

	Geheilt	Nicht Geheilt
Treatment	90	10
Kontrolle	35	35

$ARR = \frac{90}{90+10} - \frac{35}{35+35} = 0,9 - 0,5 = 0,4 \rightarrow NNT = \frac{1}{0,4} = 2,5$   
 $\Rightarrow$  Drei weitere Personen müssten behandelt werden, damit eine zusätzliche Person einen Nutzen hat (d.h. die andernfalls nicht geheilt würde).

---

# Numbers Needed to Treat (NNT)

- Auch wenn das Ziel von  $NNT$  eine einfache laienverständliche Kommunikation der Treatmenteffizienz ist, darf angezweifelt werden, ob dies immer der Fall ist, wie durch folgendes Beispiel demonstriert<sup>13</sup>:

Consider a situation in which drug versus placebo response rates are 12% versus 1%, respectively; the advantage for the drug is 11%, and the NNT is 9. Consider another situation in which the drug versus placebo response rates are 99% versus 88%, respectively; the NNT is again 9. These 2 situations are strikingly different. In the first situation, there is almost no placebo response, and medication is associated with a relatively large treatment gain. In the second situation, there is a large placebo response, and medication is associated with a relatively small treatment gain. Yet, the NNT is the same in the 2 situations. So, it is really important for clinicians to know not only what the unique contribution of the drug is (NNT) but also what the placebo response and nonresponse rates are.

# Odds Ratio (OR)

- Das **Odds Ratio** ( $OR$ ) vergleicht das *Heilerfolgsverhältnis in der Treatmentgruppe* zum *Heilerfolgsverhältnis in der Kontrollgruppe*:

	male	female	$\Sigma$
nonsmoker	a	b	a+b
smoker	c	d	c+d
$\Sigma$	a+c	b+d	n

$$OR = \frac{\frac{A}{B}}{\frac{C}{D}} = \frac{A \cdot D}{B \cdot C}$$

Beachte: Als Heilerfolgsverhältnis wird hier das Verhältnis der Zahl der geheilten Patienten gegenüber der Zahl der nicht geheilten Patienten verstanden.

- Hat das Treatment keine Auswirkung, so sind die Heilerfolgsverhältnisse in beiden Gruppen gleich, d.h.  $OR = 1$ .
- Ist das Treatment erfolgreich, ist das Heilerfolgsverhältnis in der Treatmentgruppen höher als in der Kontrollgruppe, d.h.  $OR > 1$ .
- Ist das Treatment sogar nachteilig, ist das Heilerfolgsverhältnis in der Treatmentgruppen *kleiner* als in der Kontrollgruppe, d.h.  $OR < 1$ .



Nota Bene

Als kleine Übung kann man das Odds Ratio für die beiden hypothetischen Beispiele auf der vorherigen Folie berechnen. Spoiler: das Odds Ratio ist für beide Fälle gleich ( $OR = 13.5$ )!

# Übersicht Effektmaße

Effekte	
Unterschiede	Zusammenhänge
2 unabhängige Messungen $d$	Intervalldaten $r$
2 abhängige Messungen $d$	Ordinaldaten $\rho / \tau$
mehr als 2 unabhängige Messungen $\eta^2$	Vierfeldertafel $\varphi$
mehr als 2 abhängige Messungen $\eta^2$	Vierfeldertafel für Erfolg/Risiko NNT / OR

# Referenz

Empfehlungen nach Lakens (2013)<sup>14</sup>:

## *d* Familie

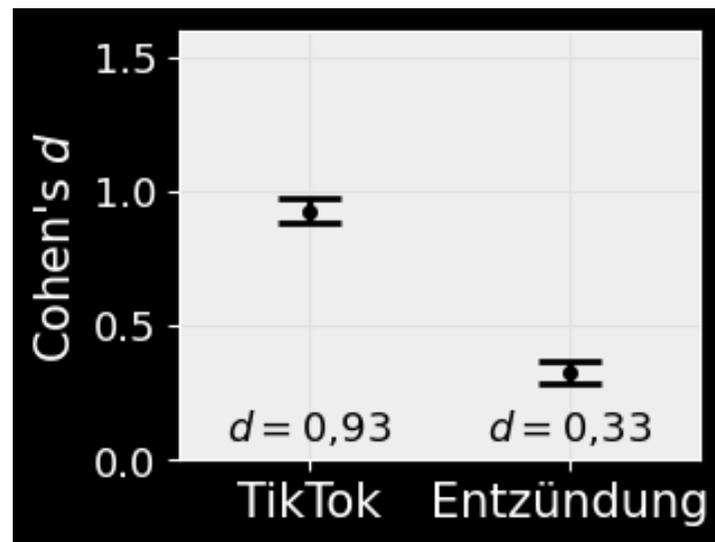
ES	Standardizer	Use
Cohen's $d_{pop}$	$\sigma$ (population)	Independent groups, use in power analyses when population $\sigma$ is known, $\sigma$ calculated with $n$
Cohen's $d_s$	Pooled $SD$	Independent groups, use in power analyses when population $\sigma$ is unknown, $\sigma$ calculated with $n-1$
Hedges' $g$	Pooled $SD$	Independent groups, corrects for bias in small samples, report for use in meta-analyses
Glass's $\Delta$	$SD$ pre measurement or control condition	Independent groups, use when experimental manipulation might affect the $SD$
Hedges' $g_{av}$	$(SD_1 + SD_2)/2$	Correlated groups, report for use in meta-analyses (generally recommended over Hedges' $g_{rm}$ )
Hedges' $g_{rm}$	$SD$ difference scores corrected for correlation	Correlated groups, report for use in meta-analyses (more conservative than Hedges' $g_{av}$ )
Cohen's $d_z$	$SD$ difference scores	Correlated groups, use in power analyses

## *r* Familie

ES (Biased)	ES (Less Biased)	Use
eta squared ( $\mu^2$ )	omega squared ( $\omega^2$ )	Use for comparisons of effects within a single study
eta squared ( $\mu_p^2$ )	omega squared ( $\omega_p^2$ )	Use in power analyses, and for comparisons of effect sizes across studies with the same experimental design.
Generalized eta squared ( $\mu_G^2$ )	Generalized omega squared ( $\omega_G^2$ )	Use in meta-analyses to compare across experimental designs



Sie berechnen nun Cohen's  $d$  für die beiden Gruppenunterschiede hinsichtlich der TikTok-Zeit und Entzündungswerte:



Hinweis: in der Abbildung wurde nicht nur die Effektstärke selbst aufgetragen, sondern auch ein Streuungsmaß der Effektstärke (Standardfehler der Effektstärke). Dazu kommen wir in Vorlesung 08.

Langsam schärft sich das Bild: während die Entzündungswerte höchstens eine mittlere Effektstärke aufweisen ( $d = 0,33$ ), ist der TikTok-Effekt beeindruckend groß:  $d = 0,93$ .

# Fußnoten

1. <https://www.laenderdaten.info/durchschnittliche-koerpergroessen.php>
2. Arbeitsberichte Dresdner Soziologie Nr. 21, <https://tud.qucosa.de/api/qucosa%3A24622/attachment/ATT-0/>
3. Cumming G (2013) Cohen's d needs to be readily interpretable: Comment on Shieh (2013). Behav Res 45:968–971.
- 4.

Delacre M, Lakens D, Ley C, Liu L, Leys C (2023) Why Hedges'  $g$ 's based on the non-pooled standard deviation should be reported with Welch's t-test. Open Science Framework. Available at: <https://osf.io/tu6mp>. Hinweis: der Terminus "ungepoolte" Varianz meint hier, dass die Daten nicht implizit in einen Pool geworfen werden und dann die Gesamtvarianz

berechnet wird; stattdessen wird unabhängig von möglicherweise unterschiedlich großen Gruppengrößen  $n_A$  und  $n_B$  der Mittelwert der beiden Einzelvarianzen berechnet.

5. Cousineau D (2020) Approximating the distribution of Cohen's  $d_p$  in within-subject designs. TQMP 16:418–421.
6. <http://jakewestfall.org/blog/index.php/2016/03/25/five-different-cohens-d-statistics-for-within-subject-designs/>
- 7.

Schäfer T, Schwarz MA (2019) The Meaningfulness of Effect Sizes in Psychological Research: Differences Between Sub-Disciplines and the Impact of Potential Biases. Frontiers in Psychology 10 Available at: <https://www.frontiersin.org/articles/10.3389/fpsyg.2019.00813>.

8. Lovakov A, Agadullina ER (2021) Empirically derived guidelines for effect size interpretation in social psychology. Eur J Soc Psychol 51:485–504.
9. Lovakov A, Agadullina ER (2021) Empirically derived guidelines for effect size interpretation in social psychology. Eur J Soc Psychol 51:485–504.
10. <https://twitter.com/ceptional/status/1687577019629142017>
11. Baguley T (2009) Standardized or simple effect size: What should be reported? British Journal of Psychology 100:603–617.
12. Mathur MB, VanderWeele TJ (2020) A Simple, Interpretable Conversion from Pearson's Correlation to Cohen's  $d$  for Continuous Exposures. Epidemiology 31:e16–e18.
- 13.

Andrade C (2015) The Numbers Needed to Treat and Harm (NNT, NNH) Statistics: What They Tell Us and What They Do Not: (Practical Psychopharmacology). J Clin Psychiatry 76:e330–e333.

Lakens D (2013) Calculating and reporting effect sizes to facilitate cumulative science: a practical primer for t-tests and ANOVAs. Frontiers in Psychology 4 Available at: <https://www.frontiersin.org/articles/10.3389/fpsyg.2013.00863>