Vorlesung Grundlagen der Diagnostik SS 25

Grundlagen der Diagnostik

Lerneinheit 4

Beobachter- und Beurteilerübereinstimmung I

© ⊕ ⊕ We are happy to share our materials openly:

The content of these <u>Open Educational Resources</u> by <u>Lehrstuhl für Psychologische Methodenlehre und Diagnostik, Ludwig-Maximilians-Universität München</u> is licensed under <u>CC BY-SA 4.0</u>. The CC Attribution-ShareAlike 4.0 International license means that you can reuse or transform the content of our materials for any purpose as long as you cite our original materials and share your derivatives under the same license.

Wahr-oder-Falsch Aussagen zur letzten Sitzung

Vorlesung Grundlagen der Diagnostik SS 25

- 1. Critical Incident Technique (CIT) basiert auf der Analyse theoretischer Modelle, ohne Rückgriff auf konkrete Situationen.
- 2. Inhaltsvalidität bedeutet, dass alle denkbaren Verhaltensweisen einer Person erfasst werden.
- 3. Beobachterdrift entsteht unter anderem durch Übungseffekte oder Ermüdung.
- 4. Die Unterscheidung zwischen Wahrnehmung, Registrierung und Beurteilung ist methodisch nicht notwendig.
- 5. Assessment Center sollten sich an den allgemeinen Standards psychologischer Diagnostik orientieren.

Vorlesung Grundlagen der Diagnostik SS 25



Vorlesung Grundlagen der Diagnostik SS 25

Rückblick und Ausblick

Objektivität und Reliabilität bei Verhaltensbeobachtung und anderen Beurteilungen (z.B. Interviews)

Statistische Analysen zur Objektivität:

• Berechnung von Maßen der Beobachter- bzw. Beurteilerübereinstimmung

Statistische Analysen zur Reliabilität:

- Wie viel Fehler passiert beim Rückschluss von den Ratings eines oder mehrerer Beobachterinnen auf das wahre Merkmal (latente Variable)?
- Beobachterinnen können sich objektiv einig sein und trotzdem ist die Messung unreliabel (z.B. Beurteilungsfehler aus LE 3).
- Bestimmte Maße der Beurteilerübereinstimmung (z.B. ICCs in LE 5) können auch als Reliabilitätsschätzung verwendet werden, aber die Trennung von Übereinstimmung und Reliabilität ist nicht immer einfach.

Vorlesung Grundlagen der Diagnostik SS 25

Indizes für Übereinstimmung abhängig vom Skalenniveau

Bei nominalskalierten Daten: Heute

- Prozentuale Übereinstimmung
- Cohens κ (kappa) und Scotts π (pi)
- Odds Ratio und Yules Y

Bei ordinal- und intervallskalierten Daten:

- Rangkorrelationen
- Intra-Klassen-Korrelation (Intra-Class-Correlation, ICC)

Hinweise:

- Wir besprechen den Fall für zwei Beobachter, aber es gibt auch Varianten für mehr als zwei.
- Für ordinalskalierte Daten gibt es eine Abwandlung von Cohens κ, das "gewichtete Cohens κ" und andere Koeffizienten (z.B. "Krippendorffs Alpha"), aber auf diese wird in der Vorlesung nicht eingegangen, da hierfür die Rangkorrelationen & ICCs besprochen werden

Vorlesung Grundlagen der Diagnostik SS 25

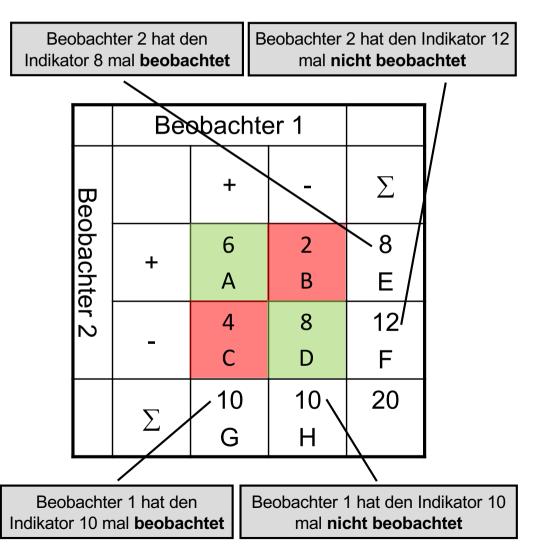
Kapitel 1: Aufbereitung der Daten

Maße für nominalskalierte Daten

Vorlesung Grundlagen der Diagnostik SS 25

Beispiel:

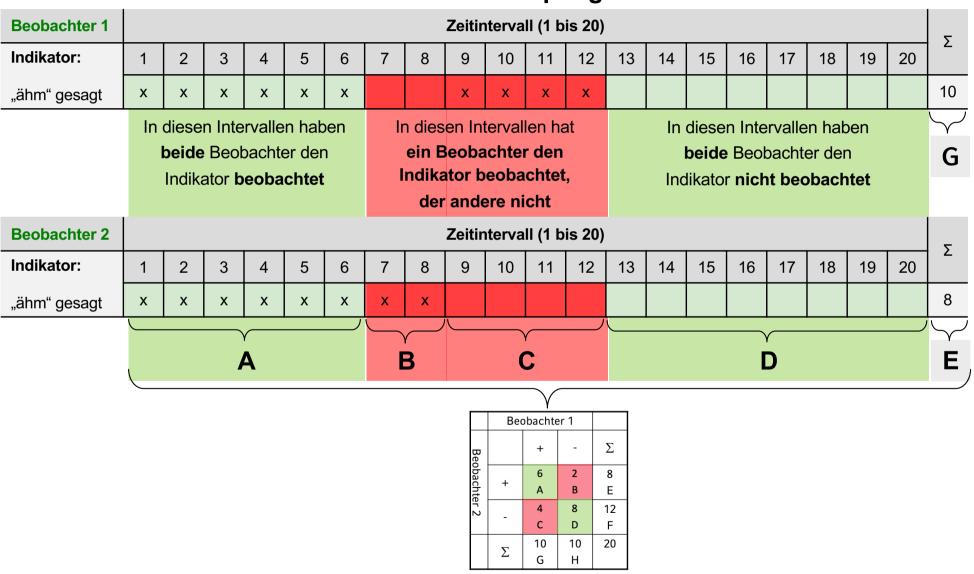
- Assessment-Center mit 20 Bewerberinnen
- Zwei Beobachterinnen (Rater)
 - → Beobachter 1 & Beobachter 2
- Urteile
 - Indikator beobachtet:
 - D.h. Kandidatin geeignet: +
 - Indikator nicht beobachtet:
 - D.h. Kandidatin nicht geeignet: -



Maße für **nominalskalierte** Daten Schritt 1: Aufbereitung der Daten

Vorlesung Grundlagen der Diagnostik SS 25

Wie kommen wir von einem Time-Sampling zu einer Kreuztabelle?



Direkte Interpretation der Kreuztabelle

Vorlesung Grundlagen der Diagnostik SS 25

- Für die Beispiele in dieser Vorlesung enthält die 2 x 2 Kreuztabelle die vollständige Information
- Wenn möglich, direkt die Kreuztabelle interpretieren, um die Qualität der Übereinstimmung zu beurteilen!

•	Wie	sehr	stimmen	die	Beobac	chter	übere	in?	١
---	-----	------	---------	-----	--------	-------	-------	-----	---

•	Wenn die Beobachter nicht übereinstimmen, welches
	Muster ergibt sich für die "Fehler"?

 Was lässt sich über die "wahre Verteilung" der zugrundeliegenden Merkmalsausprägungen vermuten?

•	Beobachten die Beobachter die
	Merkmalsausprägungen unterschiedlich häufig, bzw.
	Unterscheiden sich die Beobachter in ihrer Strenge?

 Sind sich die Beobachter hinsichtlich der Rangreihe einig (trotz unterschiedlicher Strenge), also welche Beobachtungen tendenziell "positiv" beurteilt werden?

	Bed			
В		+	-	Σ
Beobachter 2	+	6 A	2 B	8 E
chter			8	12
2	1	4 C	D	F
	Σ	10 G	10	20
	_	G	Н	N

Unterschiedliche Wahrnehmungsschwellen?

Konsistente Beurteilung?

Quantifizierung der Übereinstimmung

Vorlesung Grundlagen der Diagnostik SS 25

- In der Praxis wünscht man sich häufig eine Quantifizierung der Qualität der Beurteilerübereinstimmung über die deskriptive Betrachtung der Kreuztabelle hinaus.
- Für eine solche Quantifizierung existieren zwei Strategien:

Modellbasierte Übereinstimmung:

- Interpretation von Parametern in speziellen Messmodellen zur Beschreibung des Beobachtungs- bzw. Beurteilungsprozesses
- Methodisch anspruchsvoll (Schätzung statistischer Modelle)
- Können wir leider nur kurz ansprechen, aber nicht im Detail behandeln

Klassische Übereinstimmungsmaße und Indizes:

- Maßzahlen für Übereinstimmung basierend auf Kreuztabelle
- Zusammenfassung der Übereinstimmung in einer einzigen Zahl
- Unterschiedliche Interpretation verschiedener Maßzahlen
- Standard in der Praxis und Fokus dieser Vorlesung

Vorlesung Grundlagen der Diagnostik SS 25

Kapitel 2: Messmodelle für dichotome Ratings

Messmodelle für dichotome Ratings

Vorlesung Grundlagen der Diagnostik SS 25

Wdhl: Essentiell paralleles bzw. Essentiell τ -äquivalentes Modell für kontinuierliche Ratings

- $X_{iPerson} = \sigma_i + \theta_{Person} + \varepsilon_{iPerson}$ für alle Rater i
- Im Gegensatz zur Vorlesung Testtheorie werden Items durch Rater ersetzt
- Messgenauigkeit quantifiziert durch Reliabilität oder Standardmessfehler
- Grundlage für Maße der Beurteilerübereinstimmung bei intervallskalierten Daten → LE 5

Exkurs: Raschmodell für dichotome Ratings

•
$$P(X_{iPerson} = 1 | \theta_{Person}) = \frac{e^{\sigma_i + \theta_{Person}}}{1 + e^{\sigma_i + \theta_{Person}}}$$
 für alle Rater i

- Annahme: θ_{Person} kontinuierlich (aber $X_{iPerson}$ diskret bzw. dichotom)
- Analogie zu Statistik II: Logistische Regression mit latenter Merkmalsausprägung der Person als Prädiktor
- Testmodelle der *Item Response Theorie* → Vorlesung im Master

Latente Klassenanalyse: Messfehlerperspektive

Vorlesung Grundlagen der Diagnostik SS 25

$$\theta = 0$$

$$X = 0$$

$$\theta = 1$$

$$X = 1$$

$$X = 1$$

- Latente Variable und Itemantwort jeweils diskret mit nur zwei Ausprägungen 0 und 1
- Wahrscheinlichkeiten p_0 , p_1 , $1-p_0$, $1-p_1$ beschreiben den Antwortprozess, z.B. $p_1=P(X=1|\theta=1)$ und $1-p_1=P(X=0|\theta=1)$
- Ideale Messung:
 - p₀, p₁ nahe an 1 und 1 − p₀, 1 − p₁ nahe an 0
 → Hohe Messgenauigkeit (wenig Messfehler)
 - p_0 , p_1 , $1 p_0$, $1 p_1$ sind jeweils für alle Rater gleich \rightarrow Keine unterschiedlichen Wahrnehmungsschwellen

Latente Klassenanalyse: Übereinstimmungsperspektive

Vorlesung Grundlagen der Diagnostik SS 25

$$\theta = eindeutig 0$$

$$\theta = uneindeutig$$

$$\theta = eindeutig 1$$

$$X = 0$$

$$X = 0$$

$$X = 1$$

- Drei latente Zustände: "eindeutig 0", "eindeutig 1" und "uneindeutig"
- Eindeutige Beobachtungen werden **immer** richtig beurteilt. Die Wahrscheinlichkeiten p und 1-p beschreiben den Antwortprozess bei uneindeutigen Beobachtungen, z.B. $p = P(X = 1 | \theta = uneindeutig)$
- Ideale Messung:
 - $P(\theta = uneindeutig)$ möglichst klein
 - p entspricht dem Anteil von "eindeutig 1" unter allen "eindeutigen"
 → Kein Bias bei uneindeutigen Beobachtungen
 - p und 1 − p sind jeweils für alle Rater gleich
 → Keine unterschiedlichen Wahrnehmungsschwellen

Vorlesung Grundlagen der Diagnostik SS 25

Kapitel 3: Berechnung der klassischen Übereinstimmungsmaße

Beobachter- und Beurteilerübereinstimmung II

Vorlesung Grundlagen der Diagnostik SS 25

Unjustierte und Justierte Maße

Unjustierte Maße:

- Bestrafen unterschiedliche Strenge (Wahrnehmungsschwellen bei nominalskalierten Daten)
- Bestrafen mangelnde Konsistenz
- Konsequenz: Absolute Übereinstimmung wird bewertet

Justierte Maße:

- Bestrafen nicht unterschiedliche Strenge
- Bestrafen ausschließlich mangelnde Konsistenz
- Konsequenz: Nur Einhaltung der Rangreihe wird bewertet ("relative Übereinstimmung")

Vorlesung Grundlagen der Diagnostik SS 25

Maße für nominalskalierte Daten

- 1. Prozentuale Übereinstimmung
- 2. Cohens κ und Scotts π
- 3. Odd's Ratio und Yules Y



Maße für nominalskalierte Daten

Vorlesung Grundlagen der Diagnostik SS 25

1. Prozentuale Übereinstimmung

Berechnung:

$$P_O = \frac{\text{H\"{a}ufigkeit der \"{U}bereinstimmung}}{\text{Anzahl aller Urteile}} \cdot 100 = \frac{\text{A+D}}{\text{A+B+C+D}} \cdot 100$$

Beispiel:

$$P_0 = \frac{6+8}{6+2+4+8} = 70\%$$

	Bed			
Be		+	-	Σ
Beobachter 2		6	2	8
chte	+	6 A	2 B	Е
er 2		4	8	12
	-	4 C	D	F
	~	10	10	20
	Σ	10 G	Η	

Im Kontext von Entscheidungsfehlern:

"Relative Häufigkeit korrekter Diagnosen" bzw. "Genauigkeit" = $\frac{(TP+TN)}{N}$

Maße für nominalskalierte Daten

Vorlesung Grundlagen der Diagnostik SS 25

1. Prozentuale Übereinstimmung

Aber Achtung!

Prozentuale Übereinstimmung kann irreführend sein, wenn Merkmale sehr häufig bzw. sehr selten beobachtet werden,

z.B. sehr häufige / seltene Krankheiten (hohe / geringe *Basisrate*)

Beispiel:

$$P_0 = \frac{17+1}{17+1+1+1} = 90\%$$

	Bed	Beobachter 1				
Bec		+	ı	Σ		
Beobachter 2	т.	17	1	18		
chte	Т	Α	В	Е		
er 2		1	1	2 F		
	-	С	D	F		
	∇	18	2	20		
	Σ	G	Н			

Maße für **nominalskalierte** Daten

Vorlesung Grundlagen der Diagnostik SS 25

2. Cohens κ und Scotts π

Cohens κ und Scotts π beachten zusätzlich die "durch Zufall" erwartete Übereinstimmung:

$$\kappa, \pi = \frac{Beobachtete \ \ddot{\mathsf{U}}bereinstimunng(P_O) - Erwartete \ \ddot{\mathsf{U}}bereinstimmung\ (P_E)}{1 - Erwartete \ \ddot{\mathsf{U}}bereinstimmung\ (P_E)}$$

- Wertebereich: -1 bis +1 (in der Praxis meist > 0)
- Zähler: "Wie viel mehr Übereinstimmung als zufällig erwartbar?"
- Nenner: "Wie viel mehr Übereinstimmung wäre maximal möglich?"
 (d.h. der größtmögliche Zähler)

Maße für nominalskalierte Daten

Vorlesung Grundlagen der Diagnostik SS 25

2. Cohens κ und Scotts π

Die "durch Zufall" erwartete Übereinstimmung (P_E) wird für Cohens κ und Scotts π und unterschiedlich berechnet

- Cohens κ: Schätzung der erwarteten Übereinstimmung aus den kombinierten Randsummen
- Scotts π: Schätzung der erwarteten Übereinstimmung aus den mittleren Randsummen

Maße für **nominalskalierte** Daten

Vorlesung Grundlagen der Diagnostik SS 25

2. Cohens κ und Scotts π

Schritt 1:

Umwandlung der absoluten Häufigkeiten in relative Häufigkeiten

	Bed	bachte	er 1				Bec	bachte	er 1	
В		+	-	Σ		В		+	-	Σ
eobachter	+	6 A	2 B	8 E		Beobachter	+	.30 A	.10 B	.40 E
ter 2	-	4 C	8 D	12 F	ŕ	iter 2	-	.20 C	.40 D	.60 F
	Σ	10 G	10 H	20			Σ	.50 G	.50 H	1

Im Beispiel:

A: 6 / 20 = .30 E: 8 / 20 = .40

B: 2 / 20 = .10 F: 12 / 20 = .60

C: 4 / 20 = .20 G: 10 / 20 = .50

D: 8 / 20 = .40 H: 10 / 20 = .50

Maße für nominalskalierte Daten

Vorlesung Grundlagen der Diagnostik SS 25

2. Cohens κ und Scotts π

Schritt 2:

Die "durch Zufall" erwartete Übereinstimmung berechnen (P_E)

Wie hoch ist die Wahrscheinlichkeit, dass die Beobachter zufällig übereinstimmen, also zufällig beide "positiv" oder beide "negativ" beurteilen?

	Bed	er 1		
В		+	ı	Σ
Beobachter 2	+	.30 A	.10 B	.40 E
	ı	.20 C	.40 D	.60 F
	Σ	.50 G	.50 H	20 N

→ Addition der Wahrscheinlichkeiten

Maße für nominalskalierte Daten

Vorlesung Grundlagen der Diagnostik SS 25

2. Cohens κ und Scotts π

Schritt 2:

Die "durch Zufall" erwartete Übereinstimmung berechnen (P_E)

Wie hoch ist die Wahrscheinlichkeit, dass die Beobachter zufällig beide "positiv" beurteilen (analog für "negativ")?

	Bed	er 1		
В		+	ı	Σ
Beobachter 2	+	.30 A	.10 B	.40 E
iter 2	1	.20 C	.40 D	.60 F
	Σ	.50 G	.50 H	20 N

→ Multiplikation der Wahrscheinlichkeiten

Maße für nominalskalierte Daten

Vorlesung Grundlagen der Diagnostik SS 25

2. Cohens κ und Scotts π

Schritt 2:

Erwartete Übereinstimmung bei Zufall berechnen (P_E)

- P_E für Cohens κ : $(E \cdot G) + (F \cdot H)$
- Annahme: Jeder Beobachter hat eine eigene Wahrscheinlichkeit, eine Ausprägung zu beobachten.
- "kombinierte Randhäufigkeiten"
- Im Beispiel:

$$P_E$$
 für Cohens $\kappa = (.40 \cdot .50) + (.60 \cdot .50) = .50$

	Bed			
ш		+	-	Σ
Beobachter 2	+	.30 A	.10 B	.40 E
	1	.20 C	.40 D	.60 F
	Σ	.50 G	.50 H	20 N

Maße für nominalskalierte Daten

Vorlesung Grundlagen der Diagnostik SS 25

2. Cohens κ und Scotts π

Schritt 2:

Erwartete Übereinstimmung bei Zufall berechnen (P_E)

• P_E für Scotts
$$\pi$$
: $\left(\frac{E+G}{2}\right)^2 + \left(\frac{F+H}{2}\right)^2$

 Annahme: Beide Beobachter haben die gleiche Wahrscheinlichkeit, eine Ausprägung zu beobachten. Damit ist der Mittelwert der beiden relativen Häufigkeiten ein sinnvoller Schätzwert für diese Wahrscheinlichkeit.

	Beobachter 1				
m		+	1	Σ	
Beobachter 2	+	.30 A	.10 B	.40 E	
iter 2	1	.20 C	.40 D	.60 F	
	Σ	.50 G	.50 H	20 N	

- "mittlere Randhäufigkeiten"
- Im Beispiel: P_E für Scotts $\pi = \left(\frac{.40 + .50}{2}\right)^2 + \left(\frac{.60 + .50}{2}\right)^2 = .505$

Maße für nominalskalierte Daten

Vorlesung Grundlagen der Diagnostik SS 25

2. Cohens κ und Scotts π

Schritt 2:

Erwartete Übereinstimmung bei Zufall berechnen (P_E)

Zusammenfassung:

- P_E für Cohens κ : $(E \cdot G) + (F \cdot H)$
 - D.h. kombinierte Randhäufigkeiten

• P_E für Scotts
$$\pi$$
: $\left(\frac{E+G}{2}\right)^2 + \left(\frac{F+H}{2}\right)^2$

• D.h. mittlere Randhäufigkeiten

	Bed	bachte	er 1				
В		+	ı	Σ			
Beobachter 2	+	.30 A	.10 B	.40 E			
nter 2	-	.20 C	.40 D	.60 F			
	Σ	.50 G	.50 H	20 N			

Maße für nominalskalierte Daten

Vorlesung Grundlagen der Diagnostik SS 25

2. Cohens κ und Scotts π

Man kann **Schritt 1 und 2** auch zusammenfassen und die Umwandlung in relative Häufigkeiten in die Formeln einbauen:

• P_E für Cohens K:
$$\left(\frac{E}{N} \cdot \frac{G}{N}\right) + \left(\frac{F}{N} \cdot \frac{H}{N}\right)$$

• P_E für Scotts
$$\pi$$
: $\left(\frac{\frac{E}{N} + \frac{G}{N}}{2}\right)^2 + \left(\frac{\frac{F}{N} + \frac{H}{N}}{2}\right)^2$

	Bed			
Е		+	-	Σ
Beobachter 2	+	6 A	2 B	8 E
ıter 2	-	4 C	8 D	12 F
	Σ	10 G	10 H	20 N

Maße für nominalskalierte Daten

Vorlesung Grundlagen der Diagnostik SS 25

2. Cohens κ und Scotts π

Schritt 3:

κ bzw. π berechnen

$$\kappa, \pi = \frac{P_O - P_E}{1 - P_E}$$

Im Beispiel:

- $P_0 = .70$
- $P_E = .50 (\kappa) \text{ bzw. } .505 (\pi)$
- Cohens $\kappa = \frac{.70 .50}{1 .50} = 0.40$
- Scotts $\pi = \frac{.70 .505}{1 .505} = 0.39$

	Bed					
Beobachter 2		+	-	Σ		
	+	.30 A	.10 B	.40 E		
	-	.20 C	.40 D	.60 F		
	Σ	.50 G	.50 H	20 N		

Maße für **nominalskalierte** Daten

Vorlesung Grundlagen der Diagnostik SS 25

2. Cohens κ und Scotts π

Interpretation

- κ , π = 1: "perfekte" Übereinstimmung
- κ , π = 0: keine Übereinstimmung
- κ , π = -1: "perfekte" Nicht-Übereinstimmung
- Positive Werte: Besser als "durch Zufall" zu erwarten wäre
- Negative Werte: Schlechter als "durch Zufall" zu erwarten wäre
- -1: keine beobachteten Übereinstimmungen bei gleichzeitig maximaler Wahrscheinlichkeit zufälliger Übereinstimmungen

Maße für nominalskalierte Daten

Vorlesung Grundlagen der Diagnostik SS 25

2. Cohens κ und Scotts π

Cohens κ ist das am <u>häufigsten</u> angewandte Maß der Übereinstimmung für nominalskalierte Daten (vgl. Wirtz, 2007).

Beurteilungsgüte

Ungefähre Richtlinie für Kappa nach McHugh (2012):

- κ > .90 = fast perfekte Übereinstimmung
- κ = .80 .90 = starke Übereinstimmung
- κ = .60 .79 = moderate Übereinstimmung
- κ = .40 .59 = schwache Übereinstimmung
- κ = .21 .39 = minimale Übereinstimmung
- $\kappa = .00 .20 = keine Übereinstimmung$

Maße für **nominalskalierte** Daten

Vorlesung Grundlagen der Diagnostik SS 25

Cohens κ und Scotts π in R

data-Objekt:

_	RaterA [‡]	RaterB [‡]
1	1	1
2	1	1
3	1	1
4	1	1
5	1	1
6	1	1
7	1	0
8	1	0
9	1	0
10	1	0
11	0	1
12	0	1
13	0	0
14	0	0
15	0	0
16	0	0
17	0	0
18	0	0
19	0	0
20	0	0

Code:

```
library(irrCAC)
kappa2.table(table(data)) # cohen's kappa
scott2.table(table(data)) # scott's pi
```

Outputs:

Maße für **nominalskalierte** Daten

Vorlesung Grundlagen der Diagnostik SS 25

2. Cohens κ und Scotts π

Spezielle Eigenschaften

- Ein Wert von 1 kann nur erreicht werden, wenn beide Randverteilungen gleich sind
- Ein Wert von -1 kann nur erreicht werden, wenn alle vier Randsummen gleich sind ("symmetrische Kreuztabelle")

Cohens κ kann aber korrigiert werden, so dass es auch bei unterschiedlichen Randverteilungen/-summen maximal werden kann (-1 und +1)

Hinweis: Für den Spezialfall einer Vierfeldertafel, bei der nur eine der 4 Zellen belegt ist, sind Cohens κ und Scotts π nicht definiert.

Maße für **nominalskalierte** Daten

Vorlesung Grundlagen der Diagnostik SS 25

Ungleiche & Gleiche Randverteilungen

	Beobachter 1			
٣		+	-	\sum
eob		3	8	11
Beobachter 2	Ŧ	Α	В	Е
nter		8	1	9
2	-	С	D	F
	\sum	11	9	20
		G	Ι	Ν

	Rec			
В		+	I	\sum
Beobachter 2	4	6	2	8
ach	+	Α	В	Ш
nter		4	8	12
2	-	4 C	D	F
		10	10	20
	Σ	G	Ι	Ν

gleiche Randverteilung

ungleiche Randverteilung

Maße für nominalskalierte Daten

Vorlesung Grundlagen der Diagnostik SS 25

Ungleiche & Gleiche Randsummen

	Bec			
<u>m</u>		+	-	\sum
Beobachter 2		3	8	11
ach	Т	Α	В	Е
nter		8	1	9
N	-	С	D	F
	∇	11	9	20
	\sum	G	Ι	Ν

	Beobachter 1			
B		+	-	\sum
eob	+	2	8	10
Beobachter 2		Α	В	Е
		8	2	10
2	-	С	D	F
		10	10	20
	\sum	G	Н	N

gleiche Randsummen

Maße für nominalskalierte Daten

Vorlesung Grundlagen der Diagnostik SS 25

Zwei Fehlerquellen kann man betrachten:

1. Unterschiedliche Wahrnehmungsschwellen (auch: Strenge):

- Ab wann wird ein Verhalten als Beobachtung registriert
- Beispiel: "ähms" zählen → die eine Person registriert bereits bei "äh", die andere will ein voll ausgeprägtes "ähm" hören
- Sichtbar an ungleichen Randverteilungen (mit "McNemar"-Test könnte die Hypothese gleicher Randverteilungen in der Population geprüft werden)

2. MangeInde Konsistenz:

- Beobachterinnen nutzen grundlegend andere Definitionen oder Kriterien
- Sichtbar, wenn Übereinstimmung nicht gut ist, obwohl die Beobachterinnen insgesamt gleiche Wahrnehmungsschwellen haben / gleich streng sind (z.B. haben es beide 5 mal beobachtet, aber jeweils bei 5 verschiedenen Personen)

Maße für **nominalskalierte** Daten

Vorlesung Grundlagen der Diagnostik SS 25

	Ве	Beobachter 1				
B		+ <	-	Σ		
Beobachter 2	+	40	61	101		
ach:	Т	Α	В	Е		
ter 2		62	37	99 F		
	-	С	D	F		
	∇	102	98	200 N		
	Σ	G	Н	N		

Maße für nominalskalierte Daten

Vorlesung Grundlagen der Diagnostik SS 25

Cohens κ und Scotts π sind <u>unjustierte</u> Maße

- Das heißt, dass sie unterschiedliche Wahrnehmungsschwellen
 - = Randverteilungen der Rater "bestrafen",
 - d.h. unjustierte Koeffizienten in diesem Fall niedriger als justierte Koeffizienten
- Sind die Randverteilungen ungleich…
 - …liegen systematische Unterschiede der Rater in der Einschätzung der Grundwahrscheinlichkeiten eines Merkmals vor!
 - ...sind unterschiedliche Wahrnehmungsschwellen (z.B. Strenge) der Rater Ursache für unterschiedliche Randverteilungen
 - …fällt Cohens κ höher als Scotts π aus
 - \rightarrow Unterschiedliche Wahrnehmungsschwellen bei Scotts π stärker bestraft
 - → diese Eigenschaft kann eventuell je nach Fragestellung genutzt werden

Maße für nominalskalierte Daten

Vorlesung Grundlagen der Diagnostik SS 25

Weitere Gründe für mangelnde Übereinstimmung:

- Sind die Randverteilungen gleich (→ gleiche Strenge) ist eine mangelnde Übereinstimmung eindeutig als ein Effekt mangelnder Konsistenz interpretierbar
 - D.h. Beobachter nutzen andere Definitionen oder Kriterien
- Vorsicht: Für einen geringen Koeffizienten gibt es aber auch andere Ursachen, wie z.B. eine hohe zufällig erwartbare Übereinstimmung durch sehr geringe/hohe Basisrate

Maße für **nominalskalierte** Daten

Vorlesung Grundlagen der Diagnostik SS 25

Weitere Einflussgröße: Basisrate

Tabelle 2: Drei Beispiele für Vierfeldertafeln bei dichotomen Einschätzungen

		Beispiel A Beurteiler B		Beispiel B Beurteiler B		Beispiel C Beurteiler B				
		0	1	Σ	0	1	Σ	0	1	Σ
r A	0	74	25	99	145	18	163	94	73	167
Beurteiler	1	24	77	101	17	20	37	4	29	33
Beu	Σ	98	102	200	162	38	200	98	102	200
	Cohens ĸ		,51			,43			,24	-

Aus: Wirtz (2006). In Petermann, F. & Eid, M. (Eds.). Handbuch der Psychologischen Diagnostik (S.371).

Gleiche	Gleiche	Ungleiche
Randverteilungen	Randverteilungen	Randverteilungen
Mittlere Basisrate von 0/1-Urteilen	Geringe Basisrate von "1"-Urteilen	Basisrate unklar

Maße für nominalskalierte Daten

Vorlesung Grundlagen der Diagnostik SS 25

Wichtige Eigenschaften von Cohens κ & Scotts π

√ "Bereinigt"

Von der durch Zufall erwarteten Übereinstimmung

✓ Unjustiert

Unterschiedliche Wahrnehmungsschwellen werden bestraft

√ Abhängig von der Basisrate

- Aber anders als die prozentuale Übereinstimmung: Wenn Merkmale sehr häufig oder sehr selten vorkommen, dann ist die bei Zufall erwartete Übereinstimmung sehr hoch und Kappa / Pi werden sehr klein
- Beide Koeffizienten gibt es auch für mehr als zwei Rater oder Ratings mit mehr als zwei Kategorien ("Krippendorffs Alpha" und "Fleiss' Kappa")

Maße für nominalskalierte Daten

Vorlesung Grundlagen der Diagnostik SS 25

Zusammenhang von Cohens κ & Scotts π mit Messmodellen für dichotome Ratings

- Messfehlerperspektive: siehe Folie 12
 - Keine direkte Entsprechung, aber...
 - Scotts π entspricht (für großes N) der Intraklassenkorrelation ICC(1,1) bei kontinuierlichen Ratings → LE 5
 - Cohens κ entspricht (für großes N) der Intraklassenkorrelation ICC(2,1) bei kontinuierlichen Ratings → LE 5
- Übereinstimmungsperspektive: siehe Folie 13
 - Annahme einer idealen Messung:
 - → Kein Bias bei uneindeutigen Beobachtungen
 - → Keine unterschiedlichen Wahrnehmungsschwellen
 - Cohens κ & Scotts π entsprechen dem Anteil der eindeutig klassifizierbaren Beobachtungen an allen Beobachtungen

Maße für **nominalskalierte** Daten

Vorlesung Grundlagen der Diagnostik SS 25

3. Odds Ratio/ Yules Y

Maße für nominalskalierte Daten

Vorlesung Grundlagen der Diagnostik SS 25

Odds Ratio/ Yules Y sind justierte Maße

- Das heißt sie betrachten nur die Konsistenz von Beobachtern:
 - Wie stark stimmen Beobachter überein, ohne dass unterschiedliche Wahrnehmungsschwellen der Beobachter eine Rolle spielen?
 - Im dichotomen Fall bedeutet Konsistenz also: Im Rahmen der aufgetretenen Beobachtungshäufigkeiten überschneiden sich die Beobachtungen
- Es ist also nicht wichtig, ob Beobachter 1 das Merkmal öfter beobachtet (z.B. mehr Personen für geeignet hält) als Beobachter 2
- Odds-Ratios sollten dementsprechend dann angewandt werden, wenn unterschiedliche Wahrnehmungsschwellen keine Rolle (!) spielen sollen

Maße für nominalskalierte Daten

Vorlesung Grundlagen der Diagnostik SS 25

Odds Ratio

Schritt 1: Odds (Chance)

Odds (Übereinstimmung) =
$$\frac{\frac{A}{G}}{1 - \frac{A}{G}} = \frac{\frac{A}{G}}{\frac{C}{G}} = \frac{A}{C}$$

Lies: B2 vergibt (+), gegeben B1 hat bereits (+) vergeben

Odds (Nicht-Übereinstimmung) =
$$\frac{\frac{B}{H}}{1-\frac{B}{H}} = \frac{\frac{B}{H}}{\frac{D}{H}} = \frac{B}{D}$$

Lies: B2 vergibt (+), gegeben B1 hat bereits (-) vergeben

	Bed						
Вес		+	-	Σ			
Beobachter 2		6 A	2	8			
chte	+	Α	В	Е			
er 2		4	8	12			
	-	4 C	D	F			
	Σ	10	10	20			
	Δ	G	Н				

Exkurs

Schritt 2: Odds Ratio (Chancenverhältnis)

$$q = \frac{Odds \, (\ddot{U}bereinstimmung)}{Odds \, (Nicht - \ddot{U}bereinstimmung)} = \cdots = \frac{A \cdot D}{B \cdot C}$$

Interpretation:

- Gegeben ein Urteil, um welchen Faktor wächst die Chance für nochmal das gleiche Urteil?
- Anders ausgedrückt: Faktor, um den die Chance für ein positives (bzw. negatives) Urteil höher ist, wenn die andere Beobachterin bereits ein positives (bzw. negatives) Urteil vergeben hat
- Vorsicht: Chance ist nicht das Gleiche wie Wahrscheinlichkeit (WK), sondern WK / (1- WK)

Maße für **nominalskalierte** Daten

Vorlesung Grundlagen der Diagnostik SS 25

Odds Ratio

- Wertebereich liegt zwischen 0 und ∞
- Interpretation
 - q = 1: Übereinstimmung und Nicht-Übereinstimmung sind gleich wahrscheinlich = Zufall
 - q < 1: Systematische Nicht-Übereinstimmung
 - q > 1: Systematische Übereinstimmung

Maße für nominalskalierte Daten

Vorlesung Grundlagen der Diagnostik SS 25

Yules Y

Koeffizient dient ausschließlich der Normierung der Odds Ratio:

Yules Y =
$$\frac{\sqrt{q}-1}{\sqrt{q}+1}$$

- Liegt zwischen -1 und +1: Ermöglicht den Vergleich mit anderen Zusammenhangsmaßen, die zwischen -1 und +1 liegen
- Scotts π, Cohens κ und Yules Y sind nur bei symmetrischen Kreuztabellen gleich (= alle vier Randsummen sind gleich)

Maße für nominalskalierte Daten

Vorlesung Grundlagen der Diagnostik SS 25

3. Odds Ratio/ Yules Y

Schritt 1:

Odds Ratio berechnen

Im Beispiel:

$$q = \frac{A \cdot D}{B \cdot C} = \frac{6 \cdot 8}{2 \cdot 4} = 6$$

- → Gegeben ein Urteil, wächst die Chance für nochmal das gleiche Urteil um den Faktor 6
- → Chance für ein positives Urteil steigt um das 6fache, wenn bekannt ist, dass die andere Beobachterin auch ein positives Urteil vergeben hat

	Bed			
В		+	-	Σ
Beobachter 2	+	6 A	2 B	8 E
iter 2	-	4 C	8 D	12 F
	Σ	10 G	10 H	20

Maße für nominalskalierte Daten

Vorlesung Grundlagen der Diagnostik SS 25

Exkurs: Odds Ratio mit logistischer Regression

1. Modellgleichung aufstellen

$$P(X_{B2} = 1 | x_{B1}) = \frac{e^{\alpha + \beta \cdot x_{B1}}}{1 + e^{\alpha + \beta \cdot x_{B1}}}$$

mit $x_{B1}, x_{B2} \in \{0,1\}$

	Bed	bachte	er 1	
Bec		+	-	Σ
oba		6	2	8
Beobachter 2	+	6 A	В	E
er 2		4	8	12
	-	4 C	D	F
	Σ	10 G	10	20
		G	Н	

2. Modellparameter schätzen

$$\hat{\beta} \approx 1.79$$

 β entspricht dem logarithmierten Odds Ratio (siehe Statistik II)

3. Odds Ratio berechnen

$$q = e^{\widehat{\beta}} = e^{1.79} \approx 6$$

Maße für **nominalskalierte** Daten

Vorlesung Grundlagen der Diagnostik SS 25

3. Odds Ratio/ Yules Y

Schritt 2:

Yules Y berechnen

$$Y = \frac{\sqrt{q} - 1}{\sqrt{q} + 1}$$

Im Beispiel:

$$Y = \frac{\sqrt{6}-1}{\sqrt{6}+1} = .42$$

	Bed	Beobachter 1					
В		+	-	Σ			
Beobachter 2	+	6 A	2 B	8 E			
iter 2	1	4 C	8 D	12 F			
	Σ	10 G	10 H	20			

Maße für **nominalskalierte** Daten

Vorlesung Grundlagen der Diagnostik SS 25

Einfluss der Basisrate im Vergleich

Tabelle 2: Drei Beispiele für Vierfeldertafeln bei dichotomen Einschätzungen

		Beispiel A		Beispiel B		Beispiel C				
		Beurteiler B		Ве	Beurteiler B		Beurteiler B			
		0	1	Σ	0	1	Σ	0	1	Σ
Pr A	0	74	25	99	145	18	163	94	73	167
Beurteiler	1	24	77	101	17	20	37	4	29	33
Beu	Σ	98	102	200	162	38	200	98	102	200
	Cohens K		,51		,43			,24		
	Odds Ratio	9,50		9,48		9,34				
	Yules Y		,51			,51		,51		

Aus: Wirtz (2006). In Petermann, F. & Eid, M. (Eds.). Handbuch der Psychologischen Diagnostik (S.371).

Maße für **nominalskalierte** Daten

Vorlesung Grundlagen der Diagnostik SS 25

Wichtige Eigenschaften von Odds Ratio & Yules Y

✓ Justiert:

- Unterschiedliche Wahrnehmungsschwellen werden nicht besonders bestraft, nur Konsistenz wird berücksichtigt
- ✓ Robust gegenüber Veränderungen der Basisrate

Vorlesung Grundlagen der Diagnostik SS 25

Kapitel 4: Beispiele und Fazit

Maße für **nominalskalierte** Daten

Vorlesung Grundlagen der Diagnostik SS 25

Beispiel 1

Eine Beratungsstelle für Hochbegabung untersucht "Hochbegabte", die aufgrund sehr guter Leistungen von Lehrerinnen in der Beratungsstelle gemeldet werden.

	Bec			
D.		+	-	\sum
Beobachter 2		17	1	18
ach	+	Α	В	Е
nter		1	1	2
N	-	С	D	F
	$\overline{\nabla}$	18	2	20
	\sum	G	Н	

Maße für nominalskalierte Daten

Vorlesung Grundlagen der Diagnostik SS 25

Beispiel 1

- Prozentuale Übereinstimmung = .90
 (90%)
- Cohens κ (kappa) = .44
- Scotts π (pi) = .44
- Yules Y = .61

	Beobachter 1				
<u>Б</u>		+	_	\sum	
Beobachter 2		17	1	18	
ach	+	Α	В	E	
nter		1	1	2	
2	_	С	D	F	
	_	18	2	20	
	Σ	G	Н		

Maße für **nominalskalierte** Daten

Vorlesung Grundlagen der Diagnostik SS 25

Beispiel 2

Mit einer Verhaltensbeobachtung soll die Anzahl an Füllwörtern beobachtet werden, aber der Indikator ist nicht konkret genug formuliert, und Beobachter 2 ist etwas unkonzentrierter als Beobachter 1.

	Bec	er 1		
_ω		+	-	\sum
Beobachter 2	-	6	2	8
ach	+	Α	В	Е
nter		4	8	12
2	•	4 C	D	F
	\sum	10	10	20
	<u></u>	G	Н	N

Maße für nominalskalierte Daten

Vorlesung Grundlagen der Diagnostik SS 25

Beispiel 2

- Prozentuale Übereinstimmung = .70
- Cohens κ (kappa) = .40
- Scotts π (pi) = .39
- Yules Y = .42

	Beobachter 1				
В		+	1	\sum	
eob	-	6	2	8	
Beobachter 2	+	Α	В	Е	
nter		4	8	12	
2	ı	4 C	D	F	
	7	10	10	20	
	\sum	G	Ι	Ν	

Maße für nominalskalierte Daten

Vorlesung Grundlagen der Diagnostik SS 25

Kommentar zu Signifikanztests

- Es existieren auch Signifikanztests für κ, π (→ χ²-Test) und Y:
 - Fragestellung: Ist Übereinstimmung in der Grundgesamtheit gleich 0?
- Der berechnete Wert der Übereinstimmung ist jedoch wesentlich aussagekräftiger
 - Signifikanztests gegen den Wert 0 kann maximal eine Mindestanforderung sein und sollte ansonsten eine untergeordnete Rolle spielen (→ Manchmal hat man vielleicht sogar die Population)
- Bsp.: Cohens κ bei großer Stichprobe ab ca. 0.30 signifikant für H_0 : $\kappa = 0$, aber erst ab 0.60-0.79 als moderat zu beurteilen

Maße für nominalskalierte Daten

Vorlesung Grundlagen der Diagnostik SS 25

Fazit

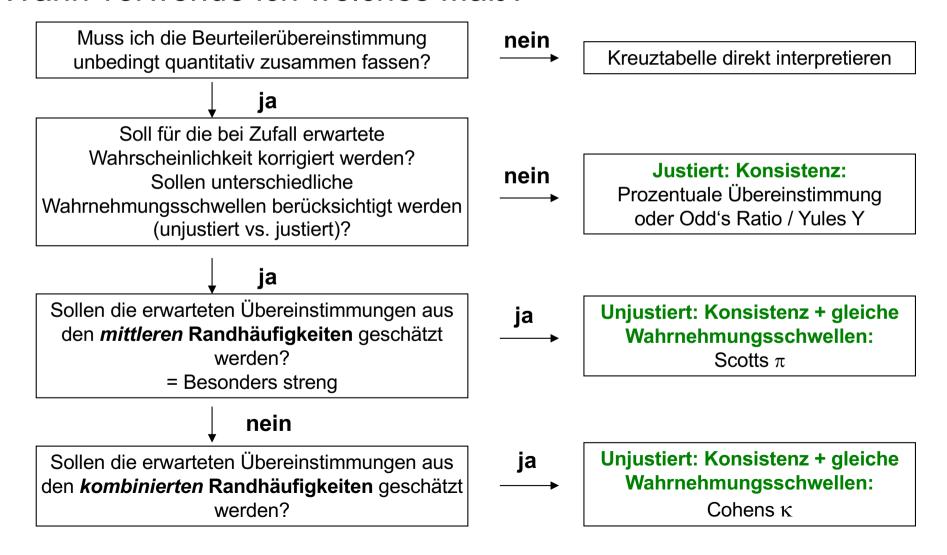
- Die heute besprochenen Maße sind auf nominalskalierte Daten anwendbar. Die Berechnungen sind mit angepassten Formeln auch für nicht-dichotome Beurteilungen möglich.
- Sollen unterschiedliche Wahrnehmungsschwellen der Rater besonders bestraft werden, dann sollte Scotts π verwendet werden, sonst Cohens κ .
- Spielt nur die Konsistenz des Ratings eine Rolle, dann sollte Yules Y verwendet werden.
- Wenn Yules Y hoch ist, aber Cohens κ / Scotts π nicht, dann weiß man dass die Konsistenz ok ist, aber die Wahrnehmungsschwellen oder die Basisrate ein Problem sind.
- Man sollte darauf achten, ob absolute oder relative Häufigkeiten vorliegen.
- Wenn möglich, sollte man immer auch direkt die Kreuztabelle betrachten und interpretieren, bevor man sich auf den komplizierten Vergleich der Übereinstimmungsmaße einlässt!

Maße für nominalskalierte Daten

Vorlesung Grundlagen der Diagnostik SS 25

Flowchart für nominalskalierte Daten

Wann verwende ich welches Maß?



Beobachter- und Beurteilerübereinstimmung I

Vorlesung Grundlagen der Diagnostik SS 25

Leitfragen zur Nachbereitung

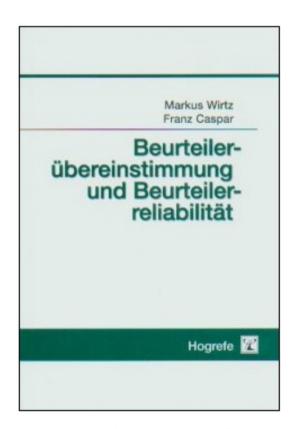
- Was ist mit Wahrnehmungsschwellen (Strenge) und Konsistenz im Rahmen der Beurteilerübereinstimmung gemeint?
- Welche klassischen Maße zur Berechnung der Beurteilerübereinstimmung für nominalskalierte Daten gibt es?
- Wie sind diese definiert? Wie werden sie berechnet?
- Wie werden sie interpretiert?
- Welche Eigenschaften haben sie?
- Wie unterscheiden sie sich voneinander? Wann nehme ich was?



Beobachter- und Beurteilerübereinstimmung I

Vorlesung Grundlagen der Diagnostik SS 25

Literatur Beobachterübereinstimmung:



Wirtz, M. & Caspar, F. (2004). Beobachterübereinstimmung (ab S. 47) Kapitel 4.1.3, 4.1.4 und Kapitel 6. Weinheim: Juventa.



Wirtz, M. & Kutschmann, M. (2006). Methoden zur Bestimmung der Beurteilerübereinstimmung. Handbuch der psychologischen Diagnostik (S. 369-380). Göttingen: Hogrefe.

https://www.researchgate.net/publication/321255287 Methoden zur B estimmung der Beurteilerubereinstimmung