

Enzyklopädie der Psychologie

Methoden der
psychologischen Diagnostik

Psychologische Diagnostik

2



Hogrefe • Verlag für Psychologie
Göttingen • Bern • Toronto • Seattle

Enzyklopädie der Psychologie

ENZYKLOPÄDIE DER PSYCHOLOGIE

In Verbindung mit der
Deutschen Gesellschaft für Psychologie

herausgegeben von

Prof. Dr. Niels Birbaumer, Tübingen
Prof. Dr. Dieter Frey, München
Prof. Dr. Julius Kuhl, Osnabrück
Prof. Dr. Wolfgang Schneider, Würzburg
Prof. Dr. Ralf Schwarzer, Berlin

Themenbereich B

Methodologie und Methoden

Serie II

Psychologische Diagnostik

Band 2

Methoden der psychologischen Diagnostik



Hogrefe • Verlag für Psychologie
Göttingen • Bern • Toronto • Seattle

Methoden der psychologischen Diagnostik

herausgegeben von

Prof. em. Dr. Lutz F. Hornke, Aachen
Prof. em. Dr. Manfred Amelang, Heidelberg
Prof. Dr. Martin Kersting, Münster



Hogrefe • Verlag für Psychologie
Göttingen • Bern • Toronto • Seattle

© 2011 Hogrefe Verlag GmbH & Co. KG
Göttingen · Bern · Wien · Paris · Oxford · Prag · Toronto
Cambridge, MA · Amsterdam · Kopenhagen · Stockholm
Rohnsweg 25, 37085 Göttingen

<http://www.hogrefe.de>

Aktuelle Informationen · Weitere Titel zum Thema · Ergänzende Materialien

Copyright-Hinweis:

Das E-Book einschließlich aller seiner Teile ist urheberrechtlich geschützt. Jede Verwertung außerhalb der engen Grenzen des Urheberrechtsgesetzes ist ohne Zustimmung des Verlags unzulässig und strafbar.

Der Nutzer verpflichtet sich, die Urheberrechte anzuerkennen und einzuhalten.

Satz: ARThür Grafik-Design & Kunst, Weimar

Format: PDF

ISBN 978-3-8409-1524-6

Nutzungsbedingungen:

Der Erwerber erhält ein einfaches und nicht übertragbares Nutzungsrecht, das ihn zum privaten Gebrauch des E-Books und all der dazugehörigen Dateien berechtigt.

Der Inhalt dieses E-Books darf von dem Kunden vorbehaltlich abweichender zwingender gesetzlicher Regeln weder inhaltlich noch redaktionell verändert werden. Insbesondere darf er Urheberrechtsvermerke, Markenzeichen, digitale Wasserzeichen und andere Rechtsvorbehalte im abgerufenen Inhalt nicht entfernen.

Der Nutzer ist nicht berechtigt, das E-Book – auch nicht auszugsweise – anderen Personen zugänglich zu machen, insbesondere es weiterzuleiten, zu verleihen oder zu vermieten.

Das entgeltliche oder unentgeltliche Einstellen des E-Books ins Internet oder in andere Netzwerke, der Weiterverkauf und/oder jede Art der Nutzung zu kommerziellen Zwecken sind nicht zulässig.

Das Anfertigen von Vervielfältigungen, das Ausdrucken oder Speichern auf anderen Wiedergabegeräten ist nur für den persönlichen Gebrauch gestattet. Dritten darf dadurch kein Zugang ermöglicht werden.

Die Übernahme des gesamten E-Books in eine eigene Print- und/oder Online-Publikation ist nicht gestattet. Die Inhalte des E-Books dürfen nur zu privaten Zwecken und nur auszugsweise kopiert werden.

Diese Bestimmungen gelten gegebenenfalls auch für zum E-Book gehörende Audiodateien.

Autorenverzeichnis

Prof. em. Dr. Manfred Amelang

Universität Heidelberg
Psychologisches Institut
Hauptstraße 47
69117 Heidelberg
E-Mail: manfred.amelang@
psychologie.uni-heidelberg.de

Prof. Dr. Martin E. Arendasy

Universität Graz
Institut für Psychologie
Universitätsplatz 2
A-8010 Graz
Österreich
E-Mail: martin.arendasy@uni-graz.at

Dr. Nicolas Becker

Universität des Saarlandes
Differentielle Psychologie und
Psychologische Diagnostik
Campus A1.3
66123 Saarbrücken
E-Mail: nicolas.becker@mx.uni-
saarland.de

Dr. Andreas Frey

Leibniz-Institut für die Pädagogik
der Naturwissenschaften (IPN)
an der Universität Kiel
Olshausenstraße 62
24098 Kiel
E-Mail: frey@ipn.uni-kiel.de

Prof. Dr. Janet Harkness

Survey Research and Methodology
Program
University of Nebraska-Lincoln
UNL Gallup Research Center
201 North 13th Street
Lincoln, NE 68588-0241
USA
E-Mail: jharkness2@unl.edu

Dr. Philipp Yorck Herzberg

Abteilung für Medizinische
Psychologie und Medizinische
Soziologie
Universitätsklinikum Leipzig AöR
Philipp-Rosenthal-Straße 55
04103 Leipzig
E-Mail: herzberg@medizin.uni-
leipzig.de

Prof. em. Dr. Lutz F. Hornke

Institut für Psychologie
RWTH Aachen University
Jägerstraße 17/19
52056 Aachen
E-Mail:
lutz.hornke@psych.rwth-aachen.de

Prof. Dr. Martin Kersting

Fachhochschule des Bundes für
öffentliche Verwaltung
Fachbereich Finanzen
Gescher Weg 100
48161 Münster
E-Mail:
Martin@Kersting-internet.de

Prof. Dr. Peter Ph. Mohler

Kirchstr. 7
76829 Leinsweiler
E-Mail: pmohler798@aol.com

Prof. Dr. Helfried Moosbrugger

Institut für Psychologie
Johann Wolfgang Goethe-Universität
Frankfurt am Main
Mertonstraße 17
60054 Frankfurt am Main
E-Mail: moosbrugger@psych.uni-
frankfurt.de

PD Dr. Beatrice Rammstedt

GESIS-Leibniz Institut
für Sozialwissenschaften
PF 122155
68072 Mannheim
E-Mail:
beatrice.rammstedt@gesis.org

Dr. Wolfgang Rauch

Institut für Psychologie
Johann Wolfgang Goethe-Universität
Frankfurt am Main
Mertonstraße 17
60054 Frankfurt am Main
E-Mail: wolfgang.rauch@psych.uni-
frankfurt.de

Prof. Dr. Jürgen Rost

Leibniz-Institut für die Pädagogik
der Naturwissenschaften (IPN)
an der Universität Kiel
Olshausenstraße 62
24098 Kiel
E-Mail: an@j-rost.de

Prof. Dr. Frank M. Spinath

Universität des Saarlandes
Differentielle Psychologie und
Psychologische Diagnostik
Campus A1.3
66123 Saarbrücken
E-Mail: f.spinath@mx.uni-
saarland.de

Mag. Markus Sommer

SCHUHFRIED GmbH
Hyrtlstraße 45
A-2340 Mödling
Österreich
E-Mail: Sommer@schuhfried.at

Jun.-Prof. Dr. Anja Strobel

TU Dresden
Professur Diagnostik und
Intervention
Helmholtzstraße 10
01069 Dresden
E-Mail: anja.strobel@tu-dresden.de

Dr. Oliver Walter

Leibniz-Institut für die Pädagogik
der Naturwissenschaften (IPN)
an der Universität Kiel
Olshausenstraße 62
24098 Kiel
E-Mail: walter@ipn.uni-kiel.de

Prof. Dr. Karl Westhoff

TU Dresden
Professur Diagnostik und
Intervention
Helmholtzstraße 10
01069 Dresden
E-Mail: Karl.Westhoff@tu-dresden.de

Dr. Safir Yousfi

Bundesagentur für Arbeit
Psychologische Forschung und
Entwicklung
Regensburger Straße 104
90478 Nürnberg
E-Mail:
Safir.Yousfi@arbeitsagentur.de

Vorwort

Eine Enzyklopädie ist zunächst sehr vereinfacht viel Arbeit, die auf viele Köpfe und fleißige Hände verteilt ist. Dabei handelt es sich um eine besondere Textsorte, die mit einer ganz besonderen fachlichen Absicht erstellt wurde.

Seit dem 5. Jh. v. Chr. versteht man unter „Enzyklopädie“ die universelle Bildung oder doch zumindest die Zusammenstellung derselben (siehe dazu den Eintrag in der wohl größten Enzyklopädie der Gegenwart, der Wikipedia, die zudem noch internetbasiert und mit vielen Querverweisen organisiert und allseits zugänglich ist). Enzyklopädisten meinten damit vom Wortstamm her zunächst einen geschlossenen Kreis an Wissen, und später stand die Idee des kreisförmig immer wiederkehrenden und zu erweiternden bzw. erneuernden Wissens im Vordergrund.

Sicher wird niemand von den vorliegenden Enzyklopädiebänden eine vollständige Darstellung von Fragen, Methodiken und Antworten der Psychologischen Diagnostik erwarten, wohl aber Hinweise auf neue Trends und Ansätze. Die Herausgeber haben versucht, dies zu ermöglichen und mit weitgefassten Vorgaben an die Autoren sowie entsprechenden Begutachtungen und Hinweisen die Grundfragen und -antworten der Psychologischen Diagnostik zu Beginn des 21. Jahrhunderts vorzustellen. Wenngleich andere, historische Enzyklopädien viel umfangreicher waren, so kann sich diese hier doch sehen lassen: 47 Artikel von 74 Autorinnen und Autoren.

Dass eine solche Arbeit mehrere Jahre in Anspruch nahm, verwundert nicht. Autoren und Herausgeber haben mit großem Elan begonnen, aber manche mussten zwischenzeitlich ihre anfänglichen Zeitplanungen angesichts ihrer auch vielfältigen anderen Aufgaben revidieren. Doch auch die traurigen und erfreulichen privaten und beruflichen Lebensereignisse verzögerten die Arbeit: Wissenschaftler, auch und gerade Psychologen, sind Menschen; Herausgeber haben das zu akzeptieren und hoffen, all dem Beruflichen und Privaten durch Benevolenz begegnet zu sein. Die Leser werden das sicher respektieren.

Jeder Beitrag ist wie bei Zeitschriften auch entsprechend begutachtet worden. Dabei war stets eine Rückmeldung anonym, und die Herausgeber haben sich entschlossen, offen ihre Korrekturen, Kommentare und Diskussionen anzubrin-

gen. Hieraus ergab sich mit einigen Autoren ein sehr fruchtbarer Dialog; die Beiträge haben dadurch fraglos gewonnen, wenngleich die Autorinnen und Autoren jeweils den eigenen Beitrag selber verantworten, denn bei ihnen liegt letztlich die Expertise, wegen derer sie zur Mitarbeit eingeladen wurden.

Im Ergebnis ist durch die Aufgeschlossenheit und die Bemühungen aller eine Enzyklopädie der Psychologischen Diagnostik entstanden, die wie ihre Vorgänger fachlich über Fragen, Methoden und Antworten informiert. Auch wenn nicht alle anfangs geplanten Beiträge geschrieben werden konnten, so liegt doch nun ein Überblick vor, der allgemein Interessierte, Studierende und Fachwissenschaftler gleichermaßen, wenngleich aus unterschiedlicher Perspektive, ansprechen kann und soll. Zumindest hat die Psychologische Diagnostik wieder eine Plattform gefunden, von der aus Gegenwarts- und Zukunftsfragen gestellt und beantwortet werden können. Die hier vorliegende Enzyklopädie zeigt die Nähe aller Autorinnen und Autoren zueinander und zum Fach. Sie ist eine Art Zugehörigkeitsausweis; denn alle, die sich hier zu Wort melden, sind die Mitgestalter und -träger zukünftiger Entwicklungen in dieser Teildisziplin der Psychologie. Dass vieles zwar auch international mitteilbar ist, aber manches doch nur im deutschsprachigen Raum rezipiert wird, das liegt an der Natur des Gebietes: Diagnostik menschlichen Erlebens und Verhaltens geht zwar in seinen allgemeinen Grundfragen über den eigenen Sprachraum hinaus, aber die gesellschaftliche Anwendung vollzieht sich nun einmal dominant innerhalb einer (Sprach-)Kultur. Beides aber bedingt sich und ist allein genommen eher problematisch.

Nach allem gebührt aus Sicht der Herausgeber allen Autorinnen und Autoren der Dank, an dieser Enzyklopädie Serie „Psychologische Diagnostik“ mitgewirkt zu haben. Die vielen Leser werden ihren Dank dadurch abtun, dass sie die Artikel zitieren, die Ideen aufgreifen, diskutieren und weiterentwickeln. Dazu gehört auch die Frage, ob die Textsorte Enzyklopädie zukünftig nicht doch auch einer anderen Präsentationsform oder einer Ergänzung durch diese bedarf.

Die Gesellschaft aber ist eingeladen, das Nachstehende zur Kenntnis zu nehmen und dort, wo es nicht unmittelbar verständlich ist, die Autorinnen und Autoren zu fragen, anzusprechen und sich mit ihnen auseinanderzusetzen. Gerade die vielen jungen Autorinnen und Autoren garantieren, dass fachlich „nachgehalten“ werden kann. Nur so wird die eigentliche Absicht einer *Enkyklios Paideia* eingelöst: Psychologie und Psychologische Diagnostik im Besonderen ist der gesellschaftlichen Wirklichkeit und Aufklärung verpflichtet; hier entfaltet sie ihre vornehmsten Wirkungen.

Lutz F. Hornke, Aachen,
Manfred Amelang, Heidelberg und
Martin Kersting, Münster

Inhaltsverzeichnis

1. Kapitel: Klassische Testtheorie. Grundlagen und Erweiterungen für heterogene Tests und Mehrfacettenmodelle

Von Wolfgang A. Rauch und Helfried Moosbrugger

1	Konzeptuelle Grundlagen und Definitionen	1
1.1	Einführung	1
1.2	Testtheorien und Voraussetzungen für deren Anwendung	2
1.3	Gütekriterien als Anforderungen an die Qualität psychologischer Diagnostik	4
1.3.1	Reliabilität als normative Anforderung	5
1.3.2	Validität als normative Anforderung	7
1.3.3	Testwert, Subtests, Testlets und Testitems	10
2	True Scores und Fehlerwerte	11
2.1	Definition der KTT	11
2.2	Einführung in das Problem der Messwiederholung	12
2.2.1	Das Uhrenbeispiel	12
2.2.2	Zur Wahrheit von „wahren“ Werten	13
2.3	Modellvorstellungen zur Entstehung der Messwertvariabilität	13
2.3.1	Zufallsziehung von Individuen	13
2.3.2	Modell der intraindividuellen Verteilung und des „stochastischen Individuums“	14
2.3.3	Zwei Dimensionen der Testwertvariabilität	15
2.4	Konstruktion des „klassischen“ Messfehlermodells	16
2.5	Experimentelle und lineare Unabhängigkeit von Testwerten	18
3	„Klassische“ Ansätze zur Schätzung von Reliabilitätskoeffizienten	20
3.1	Ausgangspunkt: Parallele Tests	20
3.2	Paralleltest-Reliabilitätsschätzung	21
3.3	Retest-Reliabilitätsschätzung	24
3.4	Reliabilitätsschätzung über Testzerlegung	24
3.4.1	Stufenweise Abschwächungen der Parallelitätsannahme	25
3.4.1.1	Parallelität	25
3.4.1.2	Tau-Äquivalenz und essenzielle Tau-Äquivalenz	25
3.4.1.3	Kongenerische Testwertvariablen	25
3.4.1.4	Nominelle Parallelität	26
3.4.2	Reliabilitätsschätzung über Testzerlegung in zwei oder drei Teile	27
3.4.2.1	Die Spearman-Brown-Formel für zwei Testhälften	27

3.4.2.2	Split-Half-Reliabilität unter (essenzieller) Tau-Äquivalenz	27
3.4.2.3	Kristofs Reliabilitätskoeffizient für drei Testteile	28
3.4.3	Reliabilitätsschätzung über Testzerlegungen in beliebig viele Teile	28
3.4.3.1	Generalisierter Spearman-Brown-Koeffizient	28
3.4.3.2	Cronbachs Alpha	29
3.4.3.3	Reliabilität bei kongenerischen Tests	29
3.4.3.4	Bias von Cronbachs Alpha bei kongenerischen Tests ..	31
3.4.3.5	Guttmans Koeffizient	32
3.4.3.6	McDonalds Koeffizient Ω	32
3.4.4	Schlussfolgerungen für die Reliabilitätsschätzung über beliebige Testzerlegungen	34
4	Erweiterung der KTT: Heterogene Tests	35
4.1	Definition von Heterogenität	35
4.2	Ursachen von Heterogenität	35
4.3	Ansätze zur Modellierung von heterogenen Tests	36
4.3.1	Zur Eindeutigkeitsproblematik bei heterogenen Testmodellen ..	36
4.3.1.1	Zur Uneindeutigkeit unterschiedlicher mehrfaktorieller Modelle	36
4.3.1.2	Zur Uneindeutigkeit von mehrfaktoriellen Modellen und Modellen mit korrelierten Fehlern	37
4.3.2	Mehrdimensionalität oder korrelierte Fehler?	38
5	Zur Rolle der konfirmatorischen Faktorenanalyse in der erweiterten KTT ..	39
5.1	Modelle mit latenten Variablen	39
5.2	Konditionale Unabhängigkeit	41
5.3	Latent-Trait-Modelle	42
5.4	Zur Unterscheidung von starken True-Score-Modellen und Latent-Trait-Modellen	42
5.5	Das konfirmatorische Faktorenmodell als Latent-Trait-Modell	43
5.5.1	Modellannahmen	43
5.5.2	Maximum-Likelihood-Schätzung und Modellfit	44
5.5.3	Relaxierte Annahmen bei konfirmatorischen Faktorenanalysen ..	44
5.5.4	Pragmatische Aspekte der konfirmatorischen Faktorenanalyse im Rahmen der erweiterten KTT	44
6	Reliabilitätsschätzung bei heterogenen Tests mit konfirmatorischen Faktorenanalysen	46
6.1	Reliabilitätsschätzung bei einfaktoriellen Modellen mit korrelierten Fehlern	47
6.1.1	„Klassische“ Reliabilität bei korrelierten Fehlern	47
6.1.2	Cronbachs Alpha bei korrelierten Fehlern	48
6.1.3	Koeffizient ω und korrelierte Fehler	48
6.2	Reliabilitätsschätzung mit mehrfaktoriellen Modellen	49
6.2.1	Das hierarchische Faktormodell	50
6.2.2	McDonalds ω und Koeffizient ω_H	52
6.2.3	Zum Einfluss der Gruppenfaktoren	52

6.3	Welcher Koeffizient sollte bei heterogenen Tests zur Reliabilitäts- schätzung gewählt werden?	55
7	Reliabilitätsschätzung bei heterogenen Tests mit explorativen Faktorenanalysen	57
7.1	Einführung	57
7.2	Die „größte untere Schranke“ der Reliabilität	58
7.3	Bentlers (2004) maximale eindimensionale Reliabilität	59
8	Stichprobentheoretische Erwägungen	60
8.1	Zur „Populationsabhängigkeit“ der KTT	60
8.2	Zu Schätzungen von Reliabilitätskoeffizienten aus Stichprobendaten ..	61
9	Mehrfacettenmodelle	62
9.1	Überblick	62
9.2	Multi-Trait-Multi-Method-Analyse	63
9.2.1	Konvergente und diskriminante Validität	63
9.2.2	Messmethoden, Methodeneffekte und Validität	64
9.2.3	Faktorenanalytische MTMM-Ansätze	66
9.3	Latent-State-Trait-Theorie	68
9.3.1	Konsistenz und (situationale) Spezifität	68
9.3.2	LST-Modell und LST-Koeffizienten	69
9.3.3	Zur Unterscheidung von Situationen und Messgelegenheiten ..	70
9.3.3.1	Bekannte oder unbekannte Situationen?	71
9.3.3.2	Mehrebenenmodelle für geschachtelte Situations- effekte	73
9.4	Generalisierbarkeitstheorie	74
9.4.1	Grundbegriffe der G-Theorie	74
9.4.2	Varianzkomponenten	75
9.4.3	Generalisierbarkeitskoeffizienten	77
9.4.4	Geschachtelte Messprozeduren	78
9.4.5	Entscheidungsstudien	79
9.4.6	Zusammenfassende Betrachtungen zur Generalisierbarkeits- theorie	79
10	Schlussfolgerungen für die Gesellschaft	80
10.1	Weiterentwicklungen der KTT	80
10.2	Konsequenzen für Lehre und Forschung	81
	Literatur	82

2. Kapitel: Psychometrische Grundlagen von Large Scale Assessments

Von Oliver Walter und Jürgen Rost

1	Testkonstruktion	88
1.1	Entwicklung der Rahmenkonzeptionen	89
1.2	Entwicklung der Testaufgaben	90

1.3	Revision der Testaufgaben, Präpilotstudien und Reviewprozesse	92
1.4	Übersetzungen der Testaufgaben	93
1.5	Signierung und Kodierung von Aufgaben mit offenem Antwortformat	94
1.6	Durchführung einer Voruntersuchung	95
2	Testdesign	96
3	Stichprobendesign	100
3.1	Definition der Zielpopulation	100
3.2	Die Entwicklung des Stichprobendesigns	102
3.3	Stichprobenfehler und Stichprobengröße	104
3.4	Ziehung der Stichprobe und Berechnung von Stichprobengewichten	107
4	Skalierung der Leistungsdaten	108
4.1	Anforderungen an die Skalierung	108
4.2	Für Large Scale Assessments geeignete Testmodelle	109
4.2.1	Logistische Testmodelle als Basismodelle in Large Scale Assessments	110
4.2.2	Verallgemeinerungen des Rasch-Modells	114
4.3	Parameterschätzung in Large Scale Assessments	118
4.3.1	Grundprinzip der Maximum-Likelihood-Methode	118
4.3.2	Itemparameterschätzung	119
4.3.3	Hintergrundmodell und latente Regression	125
4.3.4	Plausible Values und Personenparameterschätzer	130
4.3.5	Linking von Skalen verschiedener Erhebungen	133
5	Reliabilität und Modellgeltung	134
5.1	Die Schätzung der Reliabilität	134
5.2	Aspekte der Modellgeltung	138
5.2.1	Globale Modellgeltungstests	139
5.2.2	Spezielle Modellgeltungstests	140
5.2.3	Itemfitmaße	141
6	Schlussfolgerungen	144
	Literatur	145

3. Kapitel: Methoden der Item- und Skalenkonstruktion Von Safir Yousfi

1	Messtheoretische Grundlagen der Testkonstruktion	151
1.1	Grundriss der psychometrischen Testtheorie	153
1.1.1	Klassische und Probabilistische Messmodelle	155
1.1.2	Messmodelle der psychometrischen Testtheorie	156
1.2	Validität	159
2	Strategien der Testkonstruktion	164
2.1	Deduktive Methode	165
2.2	Induktive Methode	167

2.3	Externale Methode	170
2.4	Vergleich der Testkonstruktionsstrategien	171
2.5	Unterscheidungsmerkmale von psychologischen Testverfahren	172
3	Generierung von Items	175
3.1	Verhaltensstichproben, Simulationen und situative Fragen	175
3.2	Prototypenansatz	177
3.3	Lexikalischer Ansatz	178
3.4	Facettentheoretische Ansätze	179
3.5	Rationale Itemkonstruktion	180
3.6	Empfehlungen für die Itemkonstruktion bei Selbstberichtsdaten	182
4	Aggregation	185
4.1	Messung anhand von einzelnen Items	186
4.2	Aggregation durch Addition oder Mittelwertberechnung	187
4.3	Aggregation durch Linearkombination	188
4.4	Weitere statistische Methoden	189
5	Selektion von Items	190
5.1	Wissenschaftliche Aspekte der Itemselektion	190
5.2	Auswahl nach Itemkennwerten	192
5.2.1	Externe Validität	192
5.2.2	Faktorielle Validität	193
5.2.3	Interne Validität (Itemfit)	194
5.2.4	Klassische Trennschärfe- und Itemschwierigkeitskoeffizienten	195
5.2.5	Probabilistische Itemparameter	197
5.2.5.1	Trennschärfekonzepte der probabilistischen Testtheorie	197
5.2.5.2	Trennschärfe und Gütekriterien in der probabilistischen Testtheorie	198
5.2.5.3	Adaptives Testen	200
5.2.6	Inhaltliche Kriterien	201
5.3	Auswahl nach Skalenkennwerten	202
5.3.1	Algorithmen	203
5.3.2	Zielvariablen	204
5.3.3	Empirie	206
5.4	Optimal Test Design	206
6	Fazit	208
	Literatur	209

4. Kapitel: Automatisierte Itemgenerierung: Aktuelle Ansätze, Anwendungen und Forschungen Von Martin Arendasy und Markus Sommer

1	Einleitung	215
1.1	Testsicherheit	215
1.2	Bedeutung unterschiedlicher Aspekte der Validität	217

2	Automatisierte Itemgenerierung	219
2.1	Klassifikation der Ansätze zur Automatisierten Itemgenerierung	219
2.1.1	Grad der inhaltlich-theoretischen Fundierung	219
2.1.2	Grad der freien Variierbarkeit der einzelnen Bauelemente	221
2.1.3	Einbeziehung einer Qualitätssicherungskomponente	223
2.1.4	Grad der Automatisierung	225
2.2	Einordnung der unterschiedlichen Ansätze der AIG in das Klassifikationsschema	226
3	Konstruktionsphasen eines Zwei-Komponenten-Itemgenerators	232
3.1	Beschreibung des zu messenden latenten Traits	232
3.2	Ableiten der lösungsrelevanten kognitiven Prozesse und Wissensstrukturen	232
3.3	Ableitung von Radicals	234
3.4	Ableitung der funktionalen Einschränkungen	234
3.5	Formulierung und empirische Untersuchung der zu überprüfenden Meilensteine	235
4	Aktuelle Anwendung eines Zwei-Komponenten-Itemgenerators	236
4.1	Definition des zu messenden latenten Traits	236
4.2	Ableiten der lösungsrelevanten kognitiven Prozesse und Wissensstrukturen	237
4.3	Beschreibung des Aufgabenmaterials	240
4.4	Ableitung der funktionalen Einschränkungen und der Radicals	241
4.5	Beschreibung des Itemgenerators	245
4.6	Formulierung der zu überprüfenden Meilensteine und aktuelle empirische Befunde	246
4.6.1	Überprüfung der Notwendigkeit der funktionalen Einschränkungen	246
4.6.2	Überprüfung der Konstruktrepräsentation der Endlosschleifen	249
4.6.3	Weiterführende Ergebnisse zur Konstruktrepräsentation	252
4.6.4	Überprüfung der nomothetischen Spanne der Endlosschleifen	254
5	Weiterführende Forschungsfragestellungen	258
5.1	Nutzen in der Phase der Itemkonstruktion	258
5.2	Nutzen bei der Kalibrierung automatisch generierter Items	259
5.2.1	Möglichkeiten der Reduktion von Kalibrierungskosten im schemabasierten Ansatz	259
5.2.2	Möglichkeiten der Reduktion von Kalibrierungskosten im elementbasierten Ansatz	263
6	Diskussion und Ausblick	267
	Literatur	270

5. Kapitel: Kriteriumsorientierte Diagnostik

Von Philipp Yorck Herzberg und Andreas Frey

1	Einleitung	281
2	Vergleich von kriteriumsorientierten und normorientierten Tests	282
3	Konstruktion kriteriumsorientierter Tests	285
4	Modelle kriteriumsorientierter Tests	286
4.1	Klassische Testtheorie	287
4.2	Binomialmodell	287
4.3	Item-Response-Modelle	289
4.3.1	Annahmen	289
4.3.2	Modellansatz	290
4.3.3	Parameterschätzung	291
4.3.4	Schluss auf das Kriterium	291
4.3.5	Auswahl optimaler Items	296
4.3.6	Fazit	297
5	Setzen von Standards	297
5.1	Setzen von Standards durch Expertenbeurteilung	298
5.1.1	Beurteilung von Testitems	298
5.1.2	Beurteilung von Testpersonen	301
5.1.3	Expertenurteilsmethoden im Vergleich	301
5.2	Setzen von Standards als Angabe von Normen und Quoten	302
5.3	Setzen von Standards in Relation zu einem Außenkriterium	303
5.4	Beurteilung klassischer Methoden der Standardsetzung	303
5.5	Neue Ansätze zum Setzen von Standards	304
5.5.1	Ein Rahmenmodell der Standardsetzung: Generic Eclectic Method	304
5.5.2	Einbezug statistischer Methoden bei der Standardsetzung	305
6	Validität kriteriumsorientierter Tests	306
6.1	Kontentvalidität	307
6.2	Konstruktvalidität	309
6.3	Kriteriumsvalidität	311
7	Reliabilität	312
7.1	Reliabilität von Klassifikationsentscheidungen	313
7.1.1	Schwellenverlust-Indizes	313
7.1.2	Quadrierte Fehlerfunktionen	315
7.1.3	Domänenwerte	316
7.1.4	Fazit und praktische Hinweise	316
8	Fazit	317
	Literatur	318

6. Kapitel: Verhaltensbeobachtung

Von Frank M. Spinath und Nicolas Becker

1	Definition und Einteilung	326
1.1	Definition der wissenschaftlichen Verhaltensbeobachtung	326
1.2	Einteilungsgesichtspunkte	327
1.2.1	Unterteilung nach Systematik der Verhaltensbeobachtung	328
1.2.2	Unterteilung nach Segmentierung des Verhaltensstroms	328
1.2.3	Unterteilung nach Art der Abbildung	330
1.2.4	Unterteilung nach Art der Verhaltensstichprobe	335
1.2.5	Rahmen- und Durchführungsbedingungen	337
2	Beobachtungsfehler, Beobachtertraining und Gütekriterien	339
2.1	Fehlerquellen bei der Verhaltensbeobachtung	339
2.1.1	Fehler zu Lasten des Beobachtungsumfeldes	341
2.1.2	Fehler zu Lasten des Beobachtungssystems	341
2.1.3	Fehler zu Lasten des Beobachters	341
2.1.3.1	Wahrnehmungsfehler	341
2.1.3.2	Interpretationsfehler	343
2.1.3.3	Erinnerungsfehler	344
2.1.3.4	Wiedergabefehler	345
2.2	Beobachtertraining	345
2.2.1	Beobachterfehlertraining	345
2.2.2	Beobachtungsdimensionstraining	346
2.2.3	Bezugsrahmentraining	346
2.2.4	Verhaltensbeobachtungstraining	346
2.2.5	Gegenüberstellung der Effektivität der Trainingsansätze	347
2.3	Gütekriterien	347
2.3.1	Objektivität und Reliabilität	347
2.3.1.1	Verfahren für Nominaldaten	349
2.3.1.2	Verfahren für Intervalldaten	351
2.3.2	Validität	353
2.3.2.1	Inhaltsvalidität	353
2.3.2.2	Kriteriumsvalidität	354
2.3.2.3	Konstruktvalidität	354
3	Differenzielle Perspektive	356
3.1	Nutzen von Verhaltensbeobachtungen in der differentiellen Psychologie	356
3.2	Exkurs: Ökologische Validität des Dispositionsbegriffes	356
3.3	Dispositionsdiagnostik durch Verhaltensbeobachtungen	359
3.3.1	Dispositionsdiagnostik durch Beobachtungsaggregation	359
3.3.2	Dispositionsdiagnostik durch Kontrolle von Störeinflüssen	362
4	Schlussfolgerungen für die Gesellschaft	365
	Literatur	367

7. Kapitel: Interview

Von Karl Westhoff und Anja Strobel

1	Begriffsbestimmung	371
2	Einflüsse auf die Validität von Interviews	373
2.1	Strukturierungsmerkmale im Interviewprozess	374
2.1.1	Planung	374
2.1.2	Durchführung	375
2.1.3	Auswertung	376
2.2	Gestaltungsmöglichkeiten im Interviewprozess	377
2.2.1	Planung	377
2.2.2	Durchführung	378
2.2.3	Auswertung	380
2.3	Strukturierte Interviews	381
2.3.1	Klinische Psychologie	381
2.3.1.1	Das Diagnostische Interview bei psychischen Störungen (DIPS)	381
2.3.1.2	WHO – Composite International Diagnostic Interview (CIDI; WHO, 1990)	383
2.3.2	Arbeits- und Organisationspsychologie	385
2.3.2.1	Das Behavior Description Interview	385
2.3.2.2	Das Situational Interview	386
2.3.2.3	Das Multimodale Interview	386
2.3.3	Die Entscheidungsorientierte Gesprächsführung	388
2.3.4	Polizeipsychologie: Das Kognitive Interview	388
3	Die Güte des Interviews	389
3.1	Methodologische Aspekte der Evaluation	389
3.2	Psychometrische Bewertungskriterien	391
3.2.1	Objektivität	391
3.2.2	Retest-Reliabilität und Interne Konsistenz	392
3.2.3	Validität	392
3.2.3.1	Zur prädiktiven Validität	393
3.2.3.2	Zur inkrementellen Validität	395
3.2.3.3	Zur Konstruktvalidität	395
3.3	Nichtpsychometrische Bewertungskriterien – Ökonomie, Nutzen, Fairness, Akzeptanz	397
4	Der Interviewer als Schlüsselfigur im Interviewprozess	398
4.1	Interviewereinflüsse und Training des Interviewers	398
4.2	Feedback als Mittel zur Qualitätssicherung im Interviewprozess	399
5	Rechtliche und ethische Rahmenbedingungen	402
5.1	Qualitätsstandards und Grundregeln	402
5.2	Rechtliche Grundlagen	402
6	Schlussfolgerungen für die Gesellschaft	403
	Literatur	404

8. Kapitel: Psychodiagnostische Verfahren im Kulturvergleich

Von Beatrice Rammstedt, Janet Harkness und Peter Ph. Mohler

1	Einleitung	415
2	Die Vergleichbarkeit psychischer Gegebenheiten zwischen Kulturen	416
2.1	Konzeptuelle Äquivalenz	417
2.2	Phänomenologische Äquivalenz und Indikatorenäquivalenz	418
2.3	Erhebungsäquivalenz	419
2.4	Itemäquivalenz und Skalenäquivalenz	419
3	Bias und Strategien zu dessen Vermeidung	421
3.1	Konstruktbias	421
3.2	Methodenbias	422
3.3	Itembias	423
4	Methodik der kulturvergleichenden psychologischen Diagnostik	426
4.1	Perspektiven	426
4.2	Methodische Vorgehensweise	427
4.2.1	Der psychometrische Ansatz	427
4.2.2	Der quasi-experimentelle Ansatz	430
5	Entwicklung psychodiagnostischer Verfahren für kulturvergleichende Untersuchungen	432
5.1	Die Entwicklung kulturvergleichend einsetzbarer Inventare	433
5.2	Übersetzung und Kulturelle Adaptation von Inventaren	434
5.2.1	Adaptation	434
5.2.2	Übersetzung	439
5.3	Überprüfung der Angemessenheit einer Adaptation	442
5.3.1	Qualitative Strategien zur Überprüfung der Angemessenheit einer Adaptation	442
5.3.2	Quantitative Strategien zur Überprüfung der Angemessenheit einer Adaptation	444
5.4	Richtlinien für die Adaptation von Testverfahren	446
5.4.1	Richtlinien zum Kontext	446
5.4.2	Richtlinien zur Testentwicklung und Adaptation	447
5.4.3	Richtlinien zur Testanwendung	449
5.4.4	Richtlinien zur Dokumentation/Interpretation	450
6	Schlussfolgerungen	452
	Literatur	453
	Autorenregister	459
	Sachregister	471

1. Kapitel

Klassische Testtheorie. Grundlagen und Erweiterungen für heterogene Tests und Mehrfacettenmodelle

Wolfgang A. Rauch und Helfried Moosbrugger

1 Konzeptuelle Grundlagen und Definitionen

1.1 Einführung

Die Klassische Testtheorie (KTT) umfasst eine Reihe von mathematisch-statistischen Modellvorstellungen zur quantifizierenden Beschreibung der Variabilität von Testwerten; Hauptziel ist die Schätzung der Reliabilität von Testwerten. Historisch gesehen stellt die KTT einen der wesentlichen Meilensteine in der Entwicklung des heutigen empirisch-psychologischen Forschungsparadigmas dar – sie fußt auf den Beiträgen von Charles Spearman (1904a, 1904b) zur psychologischen Messung und zur Intelligenzforschung. In den letzten Jahrzehnten hätte mancher Forscher die KTT auch gerne „nur“ in den Rahmen einer historischen Abhandlung verwiesen, war doch die KTT Gegenstand z. T. vehementer Kritik und offener Ablehnung. In diesem Trend wurde die KTT auch implizit oder explizit als minderwertige Alternative zu Modellen der Item-Response-Theorie (IRT) dargestellt, vor allem deshalb, weil die KTT die Frage nicht behandelt, *wie* die Testwerte, mit denen sie sich beschäftigt, zustande kommen. Diese Kritik ist zumindest in gewisser Hinsicht nur schwer haltbar, da der wesentlichste Unterschied zwischen der „klassischen“ KTT und Modellen der IRT in der Spezifikation von Verteilungsannahmen für die Testwerte bzw. für Antworten auf Testitems liegt. Daraus folgend erlauben Modelle der IRT eine stringenter Überprüfung der Modellannahmen. Diese zusätzlichen Aspekte rechtfertigen aber nur kaum, dass diesen Modellen ein gewissermaßen qualitativ hervorgehobener wissenschaftlicher Status zugewiesen wird, wie es manche frühen Verfechter

getan haben. Allerdings gibt es einen quantitativen Vorteil von IRT-Modellen: Die stringenten Annahmen erlauben, sofern sie empirisch erfüllt sind, auch eine genauere, modellkonforme individualdiagnostische Auswertung, d. h. die Schätzung von „latenten“ Merkmalsausprägungen von Individuen.

Die Sinnhaftigkeit einer ausführlichen Abhandlung über die KTT ergibt sich also weniger hinsichtlich ihrer Anwendung für individualdiagnostische Zwecke; entsprechend werden keinerlei Ergebnisse dafür berichtet. Die hier eingenommene Perspektive ist die einer Komplementarität von KTT und IRT (vgl. auch Holland & Hoskens, 2003): Der wesentliche Nutzen der KTT und vor allem ihrer Erweiterungen im Bereich faktorenanalytischer Modelle und deren Fortführung in linearen Strukturgleichungsmodellen (SEM, *structural equation modeling*) sowie im Rahmen der Generalisierbarkeitstheorie (G-Theorie) liegt demnach in der Untersuchung von theoretischen Annahmen über die Struktur von Merkmalen in Populationen von Untersuchungsobjekten. Zumindest noch heute erlauben lineare Strukturgleichungsmodelle deutlich komplexere strukturelle Analysen als es mit IRT-Modellen bislang möglich ist. Solche Begrenzungen bei strukturellen Analysen im Rahmen von IRT-Modellen sind aber nur durch mangelnde „Computerpower“ bedingt und werden in Zukunft wohl immer weniger bestehen. Bereits heute gibt es theoretische und z. T. auch praktische Arbeiten zur Verbindung von Strukturgleichungsmodellen auf der einen Seite und generalisierten Modellen für Antwortvariablen aus der IRT auf der anderen Seite (vor allem Skrandal & Rabe-Hesketh, 2004), so dass in Zukunft die Trennung zwischen „klassischen“ Testtheorien und IRT-Modellen wohl nicht mehr sinnvoll sein wird. Vorweg geschickt sei, dass im vorliegenden Kapitel pragmatische Aspekte der Testkonstruktion nur am Rande behandelt werden. Für aktuelle Einführungen sei auf Bühner (2006) oder Moosbrugger und Kelava (2007) verwiesen. Bevor in Abschnitt 2 die KTT im engeren Sinne eingeführt wird, sollen noch kurz einige Grundlagen eingeführt und zentrale Grundbegriffe definiert werden.

1.2 Testtheorien und Voraussetzungen für deren Anwendung

Hauptanwendungsgebiet und historischer Ausgangspunkt von Testtheorien ist die psychologische Diagnostik, nämlich die

„... regelgeleitete Sammlung und Verarbeitung von gezielt erhobenen Informationen, die für das Verständnis menschlichen Verhaltens und Erlebens bedeutsam sind. Die Fragestellungen der Psychologischen Diagnostik können sich dabei auf die Beschreibung und Klassifikation, Erklärung, Vorhersage und Evaluation von Zuständen und/oder Verläufen beziehen“ (Eid & Petermann, 2006, S. 16).

Testtheorien können allerdings auch auf völlig andere Inhaltsgebiete übertragen werden als auf die „Fragestellungen der Psychologischen Diagnostik“, denn „Beschreibung und Klassifikation, Erklärung, Vorhersage und Evaluation“ charakterisieren viel allgemeiner die typischen Fragestellungen der meisten empirischen Wissenschaften. Entsprechend finden testtheoretische Ergebnisse keinesfalls nur in der psychologischen Diagnostik, sondern etwa auch in der wirtschafts- und sozialwissenschaftlichen Forschung oder der Biologie und Ökologie ihre Anwendung.

Die folgenden zwei Annahmen können als Grundvoraussetzungen für den Einsatz eines diagnostischen Verfahrens gesehen werden: (1) Personen – oder grundsätzlicher: Untersuchungsobjekte – unterscheiden sich voneinander, sonst wäre der Einsatz diagnostischer Verfahren sinnlos und (2) diese Unterschiede müssen (a) in irgendeiner Art von beobachtbaren Verhaltensäußerungen oder sonstigen empirisch erfassbaren Eigenschaften der Untersuchungsobjekte bestehen oder (b) mit den Verhaltensäußerungen oder beobachtbaren Eigenschaften zusammenhängen. Die Differenzierung von direkt oder indirekt beobachtbaren Unterschieden stellt wohl eines der zentralen Probleme psychologischer Messungen im engeren Sinne dar – gerade psychologische Tests zielen fast immer auf die Erfassung sogenannter *latenter*, d. h. nicht direkt beobachtbarer Konstrukte ab. Ein zentraler Problembereich von Testtheorien ist daher die Frage, ob die empirisch beobachteten Unterschiede zwischen den Untersuchungsobjekten auch tatsächlich in der theoretisch angenommenen Weise mit den Unterschieden auf den latenten Konstrukten zusammenhängen.

In der praktischen Anwendung werden latente Konstrukte unter Umständen „nur“ als nützliche Datenreduktion angesehen werden. In der inhaltlich psychologischen Forschung aber wird latenten Konstrukten zumindest implizit, und meist auch explizit eine objektive Realität zugeschrieben (vgl. Borsboom, Mellenberg & van Heerden, 2003). Die unterschiedlichen Konzeptionen von latenten Konstrukten werden im Verlaufe dieses Beitrags immer wieder aufgegriffen, da sie von zentraler Bedeutung nicht nur für die Bewertung unterschiedlicher testtheoretischer Ergebnisse, sondern auch für inhaltlich-psychologische Anwendungen sind.

Im engeren Sinne lässt sich eine Testtheorie als eine Reihe von mathematisch-statistischen und praxisorientierten Regeln zum Umgang mit und zur Beschreibung von Testwerten und deren Zusammenhängen untereinander definieren. Gleichzeitig umfassen Testtheorien üblicherweise auch statistische und inhaltlich-interpretative Regeln zur Beurteilung der Qualität der Testwerte bzw. der bei deren Erhebung verwendeten Prozeduren. Mit *Testwert* ist hier eine aufgrund irgendeiner Verrechnungsvorschrift zustande gekommene Zahl gemeint, die im Allgemeinen eine Zusammenfassung des Ergebnisses eines Psychologischen Tests

oder dessen einzelner Bestandteile darstellt. Grundsätzlicher kann mit dem Begriff Testwert auch das zahlenmäßige Ergebnis irgendeiner empirischen Untersuchung mit dem Ziel der Unterscheidung zwischen den Untersuchungsobjekten gemeint sein. Testtheorien beschäftigen sich also bereits mit dem Ergebnis einer Messung und sind damit von *Messtheorien* abzugrenzen, die sich mit der Frage beschäftigen, inwieweit sich Beziehungen zwischen den Objekten empirischer Forschung durch die Beziehungen zwischen Zahlen repräsentieren lassen.

Da Testtheorien aber hauptsächlich für die Analyse von Psychologischen Tests im engeren Sinne verwendet werden, soll hier ebenfalls eine Definition des Begriffes „*Psychologischer Test*“ in enger Anlehnung an Rost (2004) aufgeführt werden:

Ein Psychologischer Test ist ein spezielles psychologisches Experiment mit dem Ziel, vergleichende Aussagen über die Personen abzuleiten. Vom Versuchsleiter bewusst hergestellte unterschiedliche Reizgegebenheiten (Testaufgaben oder Items) bilden die *unabhängige Variable* (UV); die Reaktion auf die Testaufgaben oder Items unter den so hergestellten verschiedenen Versuchsbedingungen wird als *abhängige Variable* (AV) beobachtet oder gemessen.

Die aufgeführte Analogie von Tests und Experimenten (vgl. Rost, 2004) ist von großem heuristischem und konzeptuellem Nutzen; die G-Theorie (vgl. Abschnitt 9.4) etwa liefert Methoden zur Untersuchung einer systematischen Variation von Test- und Itemeigenschaften, die nicht nur, aber vor allem auch einem experimentellen Ansatz dienen. Insgesamt ist die gemeinsame Untersuchung von interindividuellen Unterschieden und Effekten experimenteller Bedingungsvariation ein häufig genanntes Ziel der Allgemeinen Psychologie wie auch der Persönlichkeitspsychologischen Forschung (vgl. Moosbrugger & Rauch, 2009). Auch aus der Perspektive der IRT wird immer mehr auf eine vergleichbare Zusammenführung der „zwei psychologischen Disziplinen“ abgezielt (Cronbach, 1957; De Boeck & Wilson, 2004).

1.3 Gütekriterien als Anforderungen an die Qualität psychologischer Diagnostik

Im vorangegangenen Abschnitt wurde als Gegenstand (und implizit auch als Zweck) von Testtheorien unter anderem die Beurteilung der Qualität von Testwerten und der Prozeduren zu ihrer Erhebung genannt. Die Qualitätsanforderungen dabei sind eine Teilmenge der grundsätzlichen Anforderungen, *wie* eine wissenschaftliche Untersuchung durchzuführen ist. Solche normativen Grundsätze entstehen üblicherweise als mehr oder weniger informeller Konsens einer

wissenschaftlichen Gemeinschaft (vgl. Messick, 1989). Gerade für die psychologische Diagnostik wurden Qualitätsgesichtspunkte in den letzten Jahren allerdings immer stärker formalisiert (vgl. Moosbrugger & Höfling, 2006). Insbesondere in der deutschsprachigen Psychologie werden solche Qualitätsanforderungen meist unter dem Oberbegriff *Gütekriterien* behandelt. Für Qualitätsanforderungen an berufsbezogene Eignungsbeurteilungen existiert sogar eine DIN-Norm (die DIN 33430; DIN, 2002).

Testtheorien haben in erster Linie eine normative Aufgabe zur Sicherung der Gütekriterien der *Reliabilität* und *Validität*. Allerdings wird die Validität in der KTT und der IRT nicht direkt thematisiert. Untersuchungen der Validität erfolgen stattdessen mit einem Methodenarsenal, das über rein testtheoretische Überlegungen hinausgeht, und dabei auf Ergebnisse aus den Testtheorien zurückgreift. Das dritte Hauptgütekriterium, die *Objektivität*, im Sinne einer möglichst weitgehenden Unabhängigkeit der Untersuchungsergebnisse von der Person desjenigen, der die Untersuchung durchführt, wird im Rahmen von Testtheorien meist implizit vorausgesetzt und nicht explizit behandelt.

Im Folgenden soll die Definition dieser Gütekriterien zunächst als normative Anforderung erfolgen, d. h. noch ohne Bezug zu einer spezifischen Testtheorie. Zu beachten ist dabei, dass der konzeptuelle Unterschied zwischen Gütekriterien als Qualitätsanforderungen und quantitativen Angaben in Form von Reliabilitäts- und Validitätskoeffizienten zur Beurteilung der Erfüllung dieser Anforderungen in der Literatur häufig nicht deutlich gemacht wird. Insbesondere der Begriff der Reliabilität ist mit der KTT verknüpft, und in der modernen, vor allem englischsprachigen psychometrischen Literatur wird der Begriff Reliabilität sehr häufig im Hinblick auf seine Formalisierung im Rahmen der KTT definiert (z. B. Lucke, 2005b; Zinbarg, Revelle, Yovel & Li, 2005).

1.3.1 Reliabilität als normative Anforderung

Zur Einführung des Reliabilitätsbegriffes wird zunächst auf zwei in der Literatur verbreitete Definitionen zurückgegriffen:

„Unter der Reliabilität eines Testes versteht man den *Grad der Genauigkeit*, mit dem er ein bestimmtes Persönlichkeits- oder Verhaltensmerkmal misst, gleichgültig, ob er dieses Merkmal auch zu messen beansprucht“ (Lienert, 1969, S. 14; Hervorhebung im Original).

„Ein Test ist dann perfekt reliabel, wenn er das Merkmal, das er misst, exakt, das heißt ohne Messfehler, misst. Dabei wird das Ausmaß der Reliabilität vom Anteil des Messfehlers an der Messung bestimmt“ (Schermelleh-Engel, Kelava & Moosbrugger, 2006, S. 421).

Gemeinsam ist beiden Definitionen, dass Reliabilität als etwas Abstufbares aufgefasst wird („Ausmaß“ bei Schermelleh-Engel et al., 2006; „Grad“ bei Lienert, 1969). Ein Unterschied liegt darin, dass in der Definition von Schermelleh-Engel et al. (2006) explizit auf das Konzept „Messfehler“ (und dessen unabänderliches Vorhandensein) rekurriert wird, während der Begriff der Genauigkeit in der Definition von Lienert (1969) prinzipiell allgemeiner ist. Ungenauigkeit kann nämlich nicht nur durch bei der Messung aufgetretene Fehler entstehen, sondern auch durch zu geringe Differenzierung: Beispielsweise kann ein Lineal zwar bei der Messung im Millimeterbereich gute Hilfen leisten, aber für deutlich kleinere Strecken meist keine sichtbare Unterscheidbarkeit herstellen. Allerdings kann der Begriff der „Genauigkeit“ auch irreführend sein, da dieser Begriff im Alltagsgebrauch nicht unbedingt auch das grundsätzliche Problem einschließt, das sich hinter dem Begriff „Messfehler“ verbirgt: *Messfehler* im engeren Sinne entstehen nämlich durch „eine Vielzahl von Faktoren, die in der Messprozedur nicht kontrolliert wurden“ (Lord & Novick, 1968, S. 38; Übersetzung durch die Autoren)¹. Immerhin macht Lienerts (1969) Reliabilitätsdefinition explizit, dass Reliabilität und Validität getrennt voneinander zu beurteilen sind („gleichgültig, ob [der Test] dieses Merkmal auch zu messen beansprucht“), was vor dem Hintergrund einer naiven, alltagssprachlichen Verbindung der Begriffe „Messfehler“ und „Zuverlässigkeit“ wichtig erscheint.

Schließlich wird ein wesentliches Ergebnis der moderneren psychometrischen Literatur in beiden Definitionen vernachlässigt: Ein Test als solcher ist nicht reliabel oder unreliabel. Die Bewertung der Reliabilität kann nur erfolgen, wenn in bestimmten Situationen erhobene Testwerte vorliegen, d. h. wenn im Sinne der obigen Definitionen situationsbedingte Reaktionen auf die Testaufgaben aufgezeichnet wurden und dabei Reaktionsunterschiede zwischen den Untersuchungsobjekten beobachtet werden können. Es ist unmittelbar ersichtlich, dass der Genauigkeitsbegriff die Unterscheidung zwischen den Reaktionen von Objekten erfordert; entsprechend muss der Aspekt der Abhängigkeit der Reliabilität von den Untersuchungsobjekten und Untersuchungssituationen berücksichtigt werden. Zusammenfassend schlagen wir die folgende Definition vor:

Reliabilität ist eine Qualitätsanforderung an Testwerte: Je weniger die Testwerte durch unkontrollierbare situative Zufallseinflüsse verfälscht werden und je feinere Unterscheidungen zwischen den Untersuchungsobjekten gemacht werden können, desto höher ist die Reliabilität.

1 Interessanterweise bezeichnet Spearman (1904b) in seinem Beitrag zur Schätzung der von „Messfehlern“ bereinigten Korrelation zweier Testreihen, der einer der wichtigsten Ausgangspunkte der psychometrischen Literatur ist, die Messfehler im heutigen Sinne als „systematic deviations“ und unterscheidet diese von „accidental deviations“, die aus dem Stichprobenfehler resultieren. Erst später setzt er den Begriff „accidental deviations“ ein, der dem heutigen Messfehlerkonzept als im experimentellen Sinne unkontrollierbare Zufallsvariation entspricht.

1.3.2 Validität als normative Anforderung

Schon Spearman (1904b) beschäftigte sich in dem bereits als Ausgangspunkt der KTT genannten Artikel mit der Frage der „Kriteriumsvalidität“, nämlich mit der Schätzung der Korrelation zweier „Testreihen“ im Zusammenhang mit der Unreliabilität der Testreihen. Dennoch gibt es Autoren, die grundsätzlich verneinen, dass die Validitätsanforderungen im Rahmen des Regelwerks der KTT überhaupt behandelt werden kann (Fischer, 1974; in gewisser Weise auch Borsboom & Mellenbergh, 2002). Unzweifelhaft ist die Operationalisierung der Reliabilitätsanforderung wohl der Kernbestandteil der KTT. Die Behandlung der Validitätsanforderung ist im Rahmen der KTT hingegen nur unter Zuhilfenahme weiterer Konzepte oder unter strengen Zusatzannahmen möglich.

Diese Herauslösung der Validitätsfrage aus dem Zusammenhang der KTT ist wohl am besten durch eine kurze historische Behandlung des Begriffs der Validität zu verstehen. Bei Spearman (1904b) und anderen wurde die aus heutiger Sicht als Teilaspekt zu verstehende Frage der „Kriteriumsvalidität“ behandelt, nämlich als Frage, inwieweit sich aus einem Testergebnis sinnvolle Vorhersagen für einen anderen Inhaltsbereich treffen lassen. Bald zeigte sich aber, dass die Beurteilung der Kriteriumsvalidität alleine höchst unbefriedigend ist. Eine Beschränkung auf die Vorhersage für andere Inhaltsbereiche vernachlässigt nämlich das Problem der theoretischen und inhaltlichen Interpretation von Testwerten. Darüber hinaus kann aber auch das Problem einer *regressio ad infinitum* auftreten, wenn das Kriterium anstelle einer theoretischen und inhaltlichen Analyse selbst wieder durch zu validierende Testwerte repräsentiert wird, etwa beim häufig genannten Beispiel der Validierung eines Intelligenztests durch Korrelation mit anderen Intelligenztests.

Aus diesen Gründen beschäftigte sich der wissenschaftliche Diskurs dann auch mit der Frage der *Inhaltsvalidität*, bei der es darum geht, (a) inwieweit ein Test inhaltlich das Gebiet abbildet, für das der Test genutzt werden soll, und/oder (b) inwieweit das im Test verlangte Verhalten das für das jeweilige Gebiet relevante Verhalten abbildet, und/oder (c) inwieweit die für die Bearbeitung des Tests notwendigen Prozesse die für das jeweilige Gebiet relevanten Prozesse repräsentieren (vgl. Messick, 1989)². Allerdings wurde auch das Konzept der Inhaltsvalidität Gegenstand harscher Kritik; beim Beispiel von Intelligenztests führt eine falsch verstandene Inhaltsvalidierung wiederum zur *regressio ad infinitum*,

2 Die Unterscheidung der drei Aspekte der Inhaltsvalidität geht in neueren Definitionen der Inhaltsvalidität meist verloren, und entsprechend wird die Inhaltsvalidität von vielen Autoren auch eher als nebensächlich betrachtet. Gerade aber die Analyse der Antwortprozesse, die für die Bearbeitung eines Tests notwendig sind, wird zuweilen als unbedingt notwendig für die Beurteilung der Validität eines Tests angesehen (Borsboom, Mellenbergh & van Heerden, 2004).

wenn unter operationalistischer und quasi-behavioristischer Perspektive „intelligentes Verhalten“ als das Verhalten definiert wird, das für die Lösung der Intelligenztestaufgaben erforderlich ist.

In den 1940er und 1950er Jahren wurde das Problem der Validität dann vor allem durch die Arbeit einer „task force“ der American Psychological Association (APA) zur Entwicklung von Teststandards (vgl. Moosbrugger & Höfling, 2006) aufgegriffen. Als Meilenstein dieser Entwicklung ist der Artikel von Cronbach und Meehl (1955) anzusehen, der den Begriff der *Konstruktvalidität* einführte. In diesem Artikel wird vor allem auch der Begriff des Konstruktes näher beleuchtet; in loser Übersetzung definieren wir:

Ein Konstrukt ist eine postulierte Eigenschaft von Untersuchungsobjekten, hinsichtlich derer sich die Untersuchungsobjekte unterscheiden³ (vgl. Cronbach & Meehl, 1955, S. 283).

Die etwas umständlich erscheinende Formulierung „eine postulierte Eigenschaft“, die von den Originalautoren übernommen wurde, kann hier auch als Vermeidung möglicher erkenntnistheoretischer Fallstricke aufgefasst werden – unabhängig von der Annahme, ob ein Konstrukt tatsächlich in einer „objektiven“ Realität existiert, oder ob es sich stattdessen etwa um eine „nützliche“ oder sonst wie geartete Konstruktion ohne Anspruch auf Entsprechung in der „realen“ physischen Welt handelt, beide Ansätze lassen sich unter dem Begriff der „Postulierung“ unterbringen. Von Bedeutung ist nur, dass sich Untersuchungsobjekte hinsichtlich des Konstruktes unterscheiden lassen.

Innerhalb des von Cronbach und Meehl (1955) gesteckten Rahmens der Konstruktvalidität genügt nicht mehr eine einzelne Korrelation, um das Ausmaß der Konstruktvalidität als Validitätskoeffizient auszudrücken. Da sich die Validitätseinschätzung auf qualitativ unterschiedliche Quellen bezieht, wird der einzelne Validitätskoeffizient durch ein Programm der Konstruktvalidierung ersetzt. Dieses Programm erscheint fast wie das grundsätzliche Forschungsprogramm der empirischen Psychologie an sich (vgl. Messick, 1989): Inwieweit lassen sich Belege und Begründungen dafür finden, dass die Testwerte auch so interpretiert werden können, wie es bei der Testkonstruktion und -anwendung intendiert war? Solche Belege und Begründungen können aus der Testbearbei-

3 Cronbach und Meehl (1955) definieren den Konstruktbezug tatsächlich direkt unter Rückgriff auf einen Test: „A construct is some postulated attribute of people, assumed to be reflected in test performance“ (S. 283); im Rahmen unserer Darstellung gehen wir allerdings davon aus, dass Konstrukte unabhängig von einem (spezifischen) Test definierbar sind. Insofern ist die Frage, inwieweit sich ein Konstrukt tatsächlich in der Testleistung abbildet, keine Frage der Definition eines Konstruktes, sondern eine Kernfrage der Testvalidierung.

tung selbst wie auch aus den Zusammenhängen mit anderen Variablen stammen, wobei diese anderen Variablen entweder als erklärende Variablen oder als abhängige (= Kriteriums-)Variablen mit den Testwerten zusammenhängen. Das heißt, die Konstruktvalidierung kann sowohl experimentelle als auch korrelative Untersuchungen umfassen und erfordert grundsätzlich die vorherige Formulierung von Hypothesen über die zu erfassenden Konstrukte und über deren Beziehungen. Die im Rahmen einer Konstruktvalidierung zu prüfenden Hypothesen können sich beispielsweise auf den Zusammenhang der Testwerte mit Testwerten in einer Kriteriumsvariablen bzw. Maßen anderer Konstrukte beziehen, aber genauso auch auf die intrapsychischen Prozesse, die für die Bearbeitung der Testaufgaben nötig sind („Inhaltsvalidität“). Dementsprechend werden in vielen modernen Darstellungen Kriteriumsvalidität und Inhaltsvalidität als untergeordnete Aspekte der Konstruktvalidität betrachtet. Die traditionelle Dreigliederung in Kriteriums-, Inhalts- und Konstruktvalidität, wie sie etwa noch bei Lienert und Raatz (1998) üblich war, kann daher mittlerweile als überholt angesehen werden (Hubleby & Zumbo, 1996), nicht nur aus der Perspektive von Messick (1989), sondern auch aus der Perspektive der aktuellen alternativen Validitätskonzeption von Borsboom, Mellenbergh und van Heerden (2004).

Der vorliegende Beitrag orientiert sich wesentlich an der Validitätskonzeption von Messick (1989), da diese im Vergleich mit anderen umfassender erscheint und sich weitestgehend von einer spezifischen erkenntnistheoretischen Position freizuhalten bemüht. In loser Anlehnung an Messick (1989) definieren wir Validität als Qualitätsanforderung daher wie folgt:

Validität ist eine Anforderung an Testwerte. Das Erheben von Testwerten dient im Allgemeinen einem gewissen Zweck; je stärker die empirischen Belege darauf hinweisen und die theoretischen Herleitungen sicherstellen, dass Interpretationen der Testwerte und Schlussfolgerungen aus diesen Interpretationen diesem Zweck gerecht werden, d. h. je besser die Interpretationen der Testwerte und Schlussfolgerungen anhand wissenschaftlicher Kriterien begründet werden können, desto höher die Validität.

Im Vergleich mit den Formulierungen von Borsboom et al. (2004; „Ein Test ist valide für eine Eigenschaft, wenn a) diese Eigenschaft existiert und b) Variation in der Eigenschaft Variation in den Testwerten kausal erzeugt“, S. 1061; Übersetzung durch die Autoren) oder mit der Definition von Lienert⁴ (1969; „Die Validität eines Testes gibt den Grad der Genauigkeit an, mit dem dieser Test dasjenige Persönlichkeitsmerkmal oder diejenige Verhaltensweise, das (die) er mes-

⁴ In Lienerts (1969) Definition ist darüber hinaus der Begriff der Genauigkeit, der ja genauso auch in der Reliabilitätsdefinition verwendet wird, äußerst unglücklich.

sen soll oder zu messen vorgibt, tatsächlich misst“, S. 16) erscheint die hier vorgenommene Definition möglicherweise zunächst kompliziert und technologisch bzw. pragmatisch-funktional, da der „Zweck“ der Erhebung von Testwerten einbezogen wird. Aber selbst derjenige Wissenschaftler, der „nur“ Erkenntnis sucht, verfolgt ein Ziel und erhebt Testwerte zu diesem Zweck. Vor allem aber wird in der hier vorgenommenen Definition wie schon bei der Definition des Konstrukt-begriffs das Problem der Existenz eines zu messenden Konstruktes ausgeklammert. Daher kann im Rahmen dieser Definition sowohl unter einer konstruktivistischen oder funktionalistischen Perspektive als auch unter einer radikal-realistischen Perspektive eine Einschätzung der Validität von Testwerten vorgenommen werden. Dennoch sei erwähnt, dass sich manche Forscher wohl am zugegebenermaßen unklaren Begriff der „Begründbarkeit“ stören mögen; die dahinter stehende Annahme ist, dass die Entscheidung darüber, was als wissenschaftlich „wahr“ oder zumindest „gesicherte Erkenntnis“ akzeptiert wird, immer auch mehr oder weniger von einer sozialen Übereinkunft der wissenschaftlichen Gemeinschaft abhängt. Die Definitionen von Lienert (1969) oder Borsboom et al. (2004) dagegen behandeln das Problem, wie die Erfüllung der jeweiligen Validitätsanforderungen gezeigt werden kann, tatsächlich überhaupt nicht. Schließlich macht die Validitätsdefinition nochmals deutlich, dass die von Hubley und Zumbo (1996) ironisch als „Dreifaltigkeit“ bezeichnete Dreiteilung von Kriteriums-, Inhalts- und Konstruktvalidität obsolet ist, da für die Begründbarkeit von Schlussfolgerungen und Handlungen aus Testwerten letztlich wissenschaftliche Belege unterschiedlicher Art gemeinsam angeführt werden können und müssen (vgl. Hartig, Frey & Jude, 2007).

Für die weitere Darstellung der KTT und ihrer Erweiterungen ist zu beachten, dass die hier gewählte Validitätskonzeption sehr viel umfassender ist als diejenige, die in psychometrischen Veröffentlichungen gerade etwas älteren Datums zu finden ist. Noch stärker als beim Begriff der Reliabilität stellen wir nicht nur an uns selbst, sondern auch an zukünftige Veröffentlichungen die Forderung, explizit und deutlich zwischen der Validität als Qualitätsanforderung und den empirischen Statistiken numerischen Größen, die als Belege für die Validität angeführt werden, zu unterscheiden.

1.3.3 Testwert, Subtests, Testlets und Testitems

Die Begriffe „Test“ und „Testwert“ wurden bereits in Abschnitt 1.2 eingeführt. Insbesondere Psychologische Tests bestehen in der Regel aus mehreren unterscheidbaren Teilen: Die kleinste Einheit, in die ein psychologischer Test zerlegt werden kann, wird im Allgemeinen als *Item* bezeichnet. Sehr häufig werden aber in einem hierarchischen Schema auch mehrere Items zu *Subtests* zusammengefasst, und zwar fast immer auf der Grundlage eines inhaltlichen Zusammen-

hangs. Beispielsweise enthielt der in PISA 2003 eingesetzte Mathematiktest drei „Subskalen“⁵ (vgl. Prenzel et al., 2004). Derselbe Test enthält aber noch eine weitere Gliederungsebene: Die „Subskalen“ sind aus mehreren (komplexen) „Aufgaben“ zusammengesetzt, deren kleinste Einheiten nun die Items darstellen, wobei sich die Items auf unterschiedliche Aspekte derselben Aufgabe beziehen. Um Konsistenz mit der angloamerikanischen Literatur zu erzielen und um Verwechslungen vorzubeugen⁶, soll für die erste Gliederungsebene oberhalb der Itemebene im Folgenden der Begriff *Testlet* verwendet werden. Im (häufigeren) Fall, dass ein Test nur eine Gliederungsebene oberhalb der Itemebene hat, werden die Begriffe „Subtest“ und „Testlet“ synonym verwendet.

Der Begriff *Testlet*, der insbesondere auch im Rahmen von IRT-Modellen populär geworden ist, wurde ursprünglich verwendet, um eine schwierigkeithomogene Itemgruppe im Rahmen eines adaptiven Tests zu bezeichnen, eine „Gruppe von Items, die sich auf ein einziges Inhaltsgebiet beziehen und als Einheit entwickelt werden“ (Wainer & Thissen, 1996, S. 190; Übersetzung durch die Autoren). In späteren Veröffentlichungen (Wainer & Thissen, 1996; Lee, Brennan & Frisbie, 2000) wurde der Umfang des Begriffes allerdings erweitert: Ein Testlet ist demnach eine Gruppe von Items, die sich auf irgendeine Art (statistisch und/oder inhaltlich) gegen Items eines anderen Testlets abgrenzen lassen. Testlets können außerdem *a priori* definiert werden, wie es etwa in der ursprünglichen Definition der Fall ist. Aber häufig zeigen empirische Ergebnisse erst *post hoc* die Notwendigkeit einer gesonderten Behandlung einzelner Testlets.

2 True Scores und Fehlerwerte

2.1 Definition der KTT

Im Speziellen definieren wir den Begriff „Klassische Testtheorie“ wie folgt:

Die Klassische Testtheorie (KTT) stellt mathematische Regeln für die Analyse der Variabilität von Testwerten zur Verfügung. Sie ist ein Bündel von Annahmen, das historisch auf unterschiedliche Weise formalisiert wurde. Die wesentlichsten Annahmen sind die folgenden: Für jeden Testwert existiert ein „wahrer“ Wert (True Score; Existenzannahme), der nicht mit dem „beobachteten“ Testwert identisch sein muss, aber mit dem beobachteten Wert linear zusammenhängt (Additivitätsannahme). Aus der Differenz von

5 Der Begriff „Skala“ oder „Subskala“ ist allerdings mehrdeutig und bezieht sich zuweilen auch auf das Ergebnis von messtheoretischen Operationen oder auf das Ergebnis der Anwendung eines IRT-Modells.

6 Die Begriffe „Item“ und „(Test-)Aufgabe“ werden – anders als im PISA-Test – meist synonym verwendet.

wahren und beobachteten Werten ergeben sich „Messfehler“, die untereinander und mit den wahren Werten unkorreliert sind.

Natürlich hat eine solche Definition den Schönheitsfehler, dass die meisten der verwendeten Begriffe wiederum definitionsbedürftig sind. Entsprechend wird ein großer Teil des folgenden Abschnitts für die zentralen Begriffe „True Score“ und „Messfehler“ verwendet. Tatsächlich wollen wir der Definition von True Scores deutlich mehr Platz einräumen als es in den meisten Darstellungen der KTT üblich ist; dieses Problem ist von fundamentaler Bedeutung für Überlegungen zur Reliabilitätsschätzung und für Fragen der Testvalidität und Generalisierbarkeit.

Darüber hinaus werden unter dem Oberbegriff KTT aber nicht nur die in der Definition genannten Konzepte von True Score und Fehlerwerten mit den daraus folgenden Ableitungen behandelt; Stumpf (1996) etwa definiert die KTT zunächst als „ein Arsenal pragmatisch orientierter Prinzipien oder Regeln zu Konstruktion, Erprobung und Evaluation psychometrischer Tests und zur Interpretation von Testergebnissen“ (S. 411). Im Rahmen eines solchermaßen erweiterten „klassischen“ Arsenal von Prinzipien zur Testkonstruktion und -evaluation werden etwa Begriffe wie die der Itemanalyse (vgl. Kelava & Moosbrugger, 2007) zumindest implizit ebenfalls dem Sachgebiet der KTT zugeordnet, auch wenn Konzepte wie „Itemschwierigkeit“, „Itemvarianz“ oder „Trennschärfe“ in Modellen der IRT meist eine sehr viel genauere Behandlung erfahren (vgl. etwa Rauch, Schweizer & Moosbrugger, 2008, für ein didaktisch orientiertes Anwendungsbeispiel).

2.2 Einführung in das Problem der Messwiederholung

Der zentrale Gesichtspunkt der KTT ist die additive Zusammensetzung von wahren Wert und Messfehler, aufbauend auf der Beobachtung, dass Testwerte bei Wiederholung der Messung variieren, selbst unter der theoretischen Annahme, dass diese stabil sein sollten. Das Grundproblem besteht in den Fragen, aus welchem Grund die (unerwünschte) Variation, nämlich die Messfehler entstehen, was die „wahren“ (nämlich nicht variierenden) Werte bedeuten und wie in diesem Zusammenhang die Reliabilität der Testwerte geschätzt werden könnte.

2.2.1 Das Uhrenbeispiel

Brennan (2001) führt in das Problem der Variabilität von Testwerten und der damit verbundenen *Unzuverlässigkeit* wie folgt ein: „Eine Person mit einer Uhr weiß, wie spät es ist; eine Person mit zwei Uhren ist nie ganz sicher“ (S. 7; Übersetzung durch die Autoren). Dieses einfach erscheinende Beispiel lässt sich heranziehen, um zentrale Konzepte und Problembereiche der KTT (aber auch der psychologischen Diagnostik im Allgemeinen) zu verdeutlichen.