

VERFAHREN ZUR EVALUATION VON SURVEY-FRAGEN: EIN ÜBERBLICK

PETER PRÜFER UND MARGRIT REXROTH

In der Umfrageforschung ist es allgemein üblich, Fragen vor ihrem endgültigen Einsatz zu testen. Viele Jahre hat man sich dabei fast ausschließlich eines Verfahrens bedient, bei dem unter Einsatz von Interviewern das Funktionieren von Fragen bei der Erhebung im Feld beobachtet wurde mit dem Ziel, eventuell vorhandene Probleme zu identifizieren. Neben diesem sogenannten Standardverfahren stehen zusätzlich eine Reihe ganz unterschiedlicher Testverfahren zur Evaluation von Fragen zur Verfügung. Es handelt sich dabei einerseits um Verfahren, deren Einsatz im Pretestbereich zwar nicht neu ist, die jedoch erst in letzter Zeit wieder aktuell wurden. Andererseits arbeitet der Umfrageforscher heute auch mit Testverfahren, die ursprünglich aus der Kognitionspsychologie stammen und heute bereits einen festen Platz in der Reihe der Verfahren zur Überprüfung von Surveyfragen eingenommen haben. Der folgende Beitrag stellt die gängigsten Verfahren vor und verweist auf deren unterschiedliche Qualitäten.

Survey research questions are usually subjected to a pretest before included in the final questionnaire. In the past, pretesting almost exclusively took place in interview situations in the field. In the last decade, more experimentally oriented techniques have been introduced based on methods and theories developed in cognitive psychology. This article presents some advantages and disadvantages of the pretest practices currently under discussion.

1. Einleitung

Wer Daten mittels Umfragen erhebt, kennt das Problem: Werden die Fragen des Fragebogens „gute“ Daten liefern, d.h. werden sie zuverlässig das messen, was sie messen sollen und damit reliable und valide Antworten liefern? Daß eine Evaluation des Fragebogens am Schreibtisch keinesfalls genügt, um gute Daten sicherzustellen, betonen beispielsweise Sudman/Bradburn (1982): „Even after years of experience, no expert can write a perfect questionnaire.“ Dafür gibt es gute Gründe: Das wesentliche Problem, mit dem sich der Fragenkonstrukteur in der Praxis konfrontiert sieht, ist der offensichtliche

Mangel an empirisch fundierten, konkreten "Konstruktionsrichtlinien". Zwar gibt es zur Konstruktion von Fragen in der methodologischen Literatur sowohl eine Reihe von ad-hoc-Regeln, Rezepten und Empfehlungen (z.B. Payne 1951; Belson 1981, 1986; Sheatsley 1983; Sudman/Bradburn 1982) als auch etliche experimentell gewonnene Ergebnisse über die Auswirkungen unterschiedlicher Frageformulierungen (vgl. z.B. Schuman/Presser 1981), sowie Ergebnisse aus dem Bereich der kognitionspsychologischen Forschung (vgl. z.B. Schwarz/Sudman 1992; Schwarz/Sudman 1994; Sudman/Bradburn/Schwarz 1996). Jedoch erweisen sich diese vorhandenen Informationen in der Praxis lediglich als hilfreich, wenn es darum geht, grobe Fehler zu vermeiden.

Daraus ergibt sich die unbefriedigende Situation, daß trotz Befolgung aller vorhandenen Regeln und Informationen bei der Konstruktion von Fragen ein Restrisiko verbleibt, das auch durch noch so große Erfahrung des Fragenkonstruktors nicht vermieden werden kann. Als Konsequenz daraus ergibt sich die Notwendigkeit, die Fragen eines Fragebogens vor deren endgültigem Einsatz einem bzw. einer Reihe von Evaluationsverfahren zu unterziehen. Sudman/Bradburn (1982) geben hierzu die Empfehlung: „If you don't have the resources to pilot test your questionnaire, don't do the study.“

Welche Verfahren zur Evaluation einzelner Fragen oder des gesamten Fragebogens stehen zur Verfügung? Mit dem vorliegenden Bericht wird erstmals ein Überblick über die gängigsten heute existierenden Verfahren zur Evaluation von Surveyfragen gegeben. Dabei werden sowohl altbekannte Techniken als auch neuere Entwicklungen vorgestellt und auf die jeweiligen Stärken und Schwächen hingewiesen.

2. Welche Verfahren stehen zur Verfügung?

Die folgende Tabelle gibt einen Überblick über die in diesem Beitrag vorgestellten Verfahren.

Tabelle 1: Verfahren zur Evaluation von Surveyfragen

Testerhebungen im Feld	Kognitive Laborverfahren	Andere Verfahren
Standard-Pretest	Think-Aloud	Focus Groups
Behaviour Coding	Probing	Experten
Problem Coding	Confidence Rating	
Random Probe	Paraphrasing	
Intensive Interview	Sorting-Verfahren	
Qualitative Interviews	Response Latency	
Analyse der Antwortverteilungen		

Split-Ballot		
--------------	--	--

2.1 Testerhebungen im Feld

2.1.1 Der Standard - Pretest

Der Begriff „Standard-Pretest“ wird in der Literatur erstmals von Oksenberg/Cannell/Kalton (1991) erwähnt. Presser/Blair (1994) verwenden den synonym zu verstehenden Begriff „Conventional Pretest“. Mitunter wird auch innerhalb der Umfrageforschung die Bezeichnung „klassischer Pretest“ oder „Beobachtungspretest“ verwendet. In jüngster Zeit findet sich in der Literatur auch der Begriff „Old-style-pretest“ (Fowler 1992).

Die Begriffe lassen zurecht vermuten, daß es sich bei dieser Methode um ein „etabliertes“, häufig angewandtes Verfahren handelt, das seit Beginn der Umfrageforschung eingesetzt wurde und allgemein als „Pretest“ bezeichnet wird. Erstaunlicherweise existieren selbst für die Durchführung eines solchen Standard-Pretests keine verbindlichen bzw. allgemein akzeptierten Regeln. In fast jedem sozialwissenschaftlichen Methodenlehrbuch findet sich zwar ein Abschnitt zum Thema „Pretest“, die Autoren geben jedoch - wenn überhaupt - höchst unterschiedliche und zum Teil auch widersprüchliche Empfehlungen bezüglich dessen Durchführung.

Im folgenden sind die unterschiedlichen Empfehlungen verschiedener Autoren bezüglich der wichtigsten Pretest-Elemente aufgeführt.

- **Stichprobe:** Die empfohlene Fallzahl, d.h. die Anzahl der zu befragenden Personen variiert von $N = 10$ bis $N = 200$ (vgl. z.B. Schrader 1971; Elliott/Christopher 1973; Friedrichs 1973; Warwick/Lininger 1975; Karmasin/Karmasin 1977; Williamson/Karp/Dalphin 1977; Wellenreuther 1982; Sheatsley 1983; Fowler 1984; Converse/Presser 1986; Schnell/Hill/Esser 1995).
- **Einsatz der Interviewer:** Während einige Autoren empfehlen, bei der Durchführung von (Standard-) Pretests ausschließlich mit erfahrenen oder speziell ausgebildeten Interviewern zu arbeiten (vgl. z.B. Atteslander 1984; Schrader 1971; Elliott/Christopher 1973; Converse/Presser 1986), plädieren andere dafür, einen „Querschnitt“ aller bei der Haupterhebung beteiligten Interviewer einzusetzen (vgl. z.B. Karmasin/Karmasin 1977; DeMaio 1983). Reine Experten als Pretestinterviewer werden z.B. von Noelle (1971) und Kidder (1981) empfohlen. Daß die Forscher bzw. die Mitglieder von Projektgruppen auch selbst Pretest-Interviews durchführen sollen, findet in der Literatur breite Zustimmung (vgl. z.B. Oppenheim 1966; Williamson/ Karp/Dalphin 1977; Sudman/Bradburn 1982; Porst 1985).

- **Informiertheit der Befragten:** Einerseits besteht die Möglichkeit, Befragte über den Testcharakter der Befragung zu informieren, andererseits kann ein (Standard-) Pretest unter den Bedingungen der Hauptstudie durchgeführt werden, ohne die Befragten über den Testcharakter in Kenntnis zu setzen. Converse/Presser (1986) führen für die beiden Varianten die Begriffe „participating pretest“ und „undeclared pretest“ ein.
- **Informationsbeschaffung:** Es existieren unterschiedliche Möglichkeiten darüber, wie der Interviewer seine jeweiligen Pretest-„Erkenntnisse“ an den Forscher weiterleitet bzw. meldet: Report des Interviewers, und zwar entweder schriftlich in Form eines sogenannten „Erfahrungsberichts“ bzw. „Pretest-Reports“ (meist über jedes durchgeführte Interview) oder mündlich als sog. „Debriefing“, und zwar entweder in Einzel-Sitzungen oder alle am Standard-Pretest beteiligten Interviewer berichten in einer gemeinsamen Sitzung über ihre Interview-Erfahrungen. Grundsätzlich besteht weiterhin die Möglichkeit, auch Befragte einem Debriefing zu unterziehen (vgl. z.B. DeMaio 1996).

Trotz fehlender empirisch fundierter Regeln gibt es zumindest eine Art Übereinstimmung darüber, wie das Grundgerüst eines Standard-Pretests beschaffen ist bzw. sein sollte. Danach zeichnet sich ein Standard-Pretest durch folgende Merkmale aus:

- Einmalige Erhebung eines Fragebogens unter möglichst realistischen Hauptstudie - Bedingungen,
- Durchführung von 20 bis 50 Interviews (Quota oder Random).
- Interviewer haben die Aufgabe, Probleme und Auffälligkeiten bei der Durchführung der Interviews zu beobachten und zu berichten.
- In der Regel handelt es sich um ein passives Verfahren, d.h. der Interviewer beobachtet nur (deshalb auch „Beobachtungspretest“), ohne aktiv zu hinterfragen.

Bei dieser Vorgehensweise liegt die Strategie bzw. das Prinzip zugrunde, aus der Reaktion bzw. Antwort der Befragten Rückschlüsse auf das Fragenverständnis zu ziehen. Dabei werden ganz allgemein Probleme der Befragten bei bzw. mit einer Frage auf Konstruktionsmängel der jeweiligen Frage zurückgeführt, während Fragen, bei denen Befragte formal „korrekt“ antworten, als „gut konstruiert“ angesehen werden.

Wiegt man Stärken und Schwächen des Standard-Pretests gegeneinander ab, ergibt sich folgendes Fazit.

Stärken des Verfahrens

Ein Standard-Pretest ist in der Regel relativ schnell und problemlos durchführbar. Der organisatorische Aufwand ist eher niedrig. Dies ist besonders dann der Fall, wenn die Befragten nach einem Quotenverfahren ausgewählt werden.

- Die Kosten sind relativ niedrig. Auch hier gilt: Der Einsatz des Quotenverfahrens wirkt sich kostenmindernd aus.
- Eine annähernd realistische Schätzung der Befragungsdauer ist möglich.

Schwächen des Verfahrens

- Der diesem Pretest-Prinzip zugrundeliegende Schluß, Fragen, die Befragte formal „korrekt“ beantworten, könnten als „gut konstruiert“ angesehen werden, ist grundsätzlich unzulässig. Belson (1981, 1986) kann beispielsweise nachweisen, daß trotz - formal - korrekter Antwort ein falsches Fragenverständnis zu Grunde liegen kann.
- Die Instruktion an die Interviewer, was sie beobachten und berichten sollen, ist meist wenig präzise. Das Produkt - die Berichte - sind dementsprechend auch sowohl inhaltlich als auch formal wenig systematisch.
- Interviewer berichten trotz intensiver Schulung bei weitem nicht alle im Standard-Pretest aufgetretenen „Probleme“ (Kreiselmaier/Prüfer/Rexroth 1989).
- Insgesamt gesehen handelt es sich beim Standard-Pretest - allein schon wegen der geringen Fallzahl - um ein sehr „grobes“ Verfahren.

2.1.2 Behaviour Coding

Dieses Evaluationsverfahren hat in den letzten Jahren einen wichtigen Platz eingenommen. Cannell et al. (1989) bezeichnen das Behaviour Coding als sinnvolle und nützliche Technik zur Identifizierung von Fragenmängeln. Das grundlegende Prinzip dieser Technik basiert auf der Klassifizierung von Verhalten. Mit Hilfe eines mehr oder weniger detaillierten Codesystems wird das Verhalten von befragter Person und Interviewer bewertet und analysiert. Ursprünglich wurde diese Technik eingesetzt, um ausschließlich Interviewerverhalten zu klassifizieren und zu bewerten. In späteren Arbeiten wurde dann auch das Befragtenverhalten vercodet (vgl. hierzu z.B. Morton-Williams/Sykes 1984; Prüfer/Rexroth 1985; Oksenberg/Cannell/Kalton 1991). Die von den einzelnen Forschern konzipierten Codesysteme zur Bewertung der Verhaltensweisen bzw. der verbalen Aktivitäten im Interview unterscheiden sich etwas in Aufbau und Detailliertheit, für alle gilt jedoch das grundsätzliche Prinzip: Das Verhalten von befragter Person und Interviewer wird mit Hilfe des Codesystems systematisch registriert. Durch diese Eigenschaft der Technik lassen sich Rückschlüsse auf die Qualität einer Frage ziehen. Damit ist die Behaviour-Coding-Technik eine ernstzunehmende Alternative, wenn es im Pretest darum geht, zur Qualitätsbestimmung

von Fragen und Instrumenten nutzbringende Techniken einzusetzen. Im folgenden wollen wir kurz darauf eingehen, wie diese Methode arbeitet und was sie leistet.

Traditionelle Vorgehensweise beim Einsatz der Behaviour-Coding-Technik

Bei dem vornehmlich in der englischsprachigen Literatur beschriebenen Behaviour Coding bewerten sogenannte „Coder“ das auf Tonband aufgezeichnete Interview, d.h. sie bewerten Interviewer - und Befragtenverhalten mittels eines Codesystems, das mehr oder weniger umfangreich sein kann, und damit mehr oder weniger differenziert Verhalten erfaßt. Als Beispiel sei ein Codesystem aufgeführt, das in der Studie von Oksenberg/Cannell/Kalton (1991) verwendet wurde. Dieses Schema sieht drei Kategorien zur Bewertung von Interviewerverhalten beim Vorlesen des Fragetextes vor (Codes E, S, M) und sieben Kategorien zur Erfassung von Befragtenverhalten (Codeziffern 1 bis 7):

Codesystem für Behaviour Coding (in Klammern die jeweiligen Bezeichnungen im Original)

Code	Verhaltensbeschreibung	
	Interviewer	
E	(Exact)	Interviewer liest Frage exakt
S	(Slight change)	Interviewer nimmt leichte Veränderungen vor
M	(Major change)	Interviewer nimmt starke Veränderungen vor
	Befragte(r)	
1	(Interruption)	Befragte(r) antwortet vorzeitig
2	(Clarification)	Befragte(r) will Wiederholung der Frage oder Klärung der Frage oder macht Bemerkung, die auf Verständnisproblem schließen läßt
3	(Adequate answer)	Befragte(r) antwortet adäquat
4	(Qualified answer)	Antwort ist adäquat, zusätzliche Bemerkung läßt jedoch auf Unsicherheit schließen
5	(Inadequate answer)	Inadäquate Antwort
6	(Don't know)	Weiß nicht
7	(Refusal to answer)	Befragte(r) verweigert Beantwortung der Frage

Vor der Vercodung sollte genau festgelegt werden, welches Verhalten überhaupt berücksichtigt werden soll bzw. welche Vercodungsregeln zugrunde gelegt werden. Oksenberg/Cannell/Kalton (1991) berichten von drei möglichen Varianten:

1. Vercodet wird nur die erste Reaktion des/der Befragten nach der Präsentation des Fragestimulus.
2. Vercodet wird das gesamte Befragtenverhalten, d.h. es können mehrere Codes pro Frage - auch mehrfach - vergeben werden.
3. Vercodet wird das gesamte Befragtenverhalten, im Unterschied zur zweiten Variante werden identische Codes nur einmal vergeben, auch wenn sie mehrmals auftreten sollten.

In der Studie von Oksenberg/Cannell/Kalton (1991) kam die dritte Variante zur Anwendung. Praktisch erhält man damit pro Frage eine Häufigkeitsverteilung über alle bei der Frage vergebenen Codes über alle Fälle. Dabei werden sowohl Art als auch Häufigkeit der durch die Codes repräsentierten Verhaltensweisen von Interviewern und Befragten als Qualitätsindikator dieser Frage gewertet.

Vergleicht man die Leistungsfähigkeit der Technik mit der anderer Pretesttechniken, so ist sie mit Abstand diejenige Technik, deren Pretesterkenntnisse am reliabelsten sind (vgl. hierzu z.B. Presser/Blair 1994). Gerade im Vergleich zum Standard-Pretest, bei dem die Pretesterkenntnisse oftmals von der subjektiven Wahrnehmung des einzelnen Interviewers geprägt sind, besticht die Behaviour-Coding-Technik durch ihre objektive und systematische Vorgehensweise. Der wesentliche Nachteil der Technik liegt allerdings darin, daß Hinweise auf mögliche Ursachen für inadäquates Verhalten nicht erfaßt werden. Dies ist umso schwerwiegender, als es ja gerade das Ziel einer Evaluation ist, ganz konkret die Schwächen bei einer Frage zu erkennen, um sie dann für den Hauptfragebogen zu eliminieren. Da das Erhebungsverfahren beim Behaviour Coding im Grunde demjenigen des Standard-Pretests entspricht (d.h. aus der Beobachtung des Befragtenverhaltens werden Rückschlüsse auf das Fragenverständnis gezogen), ist ein weiterer Nachteil, daß trotz formal korrekter Antwort ein falsches Fragenverständnis zu Grunde liegen kann. Schließlich besteht ein Nachteil des Behaviour Coding darin, daß die Interviews auf Band aufgezeichnet werden müssen.

Eine Variante des Behaviour Coding: Problem Coding

In der Feldabteilung von ZUMA wird das Behaviour Coding seit Jahren modifiziert eingesetzt. Die Modifikation besteht in einer Verbindung zwischen Standard-Pretest und traditionellem Behaviour Coding. Die Autoren nennen dieses Verfahren „Problem Coding“. Wesentlich für das Problem Coding ist, daß die Bewertung der Verhaltensweisen des/der Befragten nicht vom Coder nach dem Interview, sondern vom Interviewer selbst

während des Interviews vorgenommen wird. Dabei ist für den Interviewer die Anwendung eines ausführlichen, detaillierten Codesystems nicht möglich. Es würde den Interviewer in Verbindung mit seinen eigentlichen Aufgaben, nämlich der korrekten Durchführung des Interviews und der Registrierung der Antworten unter Berücksichtigung der eingeübten Regeln zur Durchführung des standardisierten Interviews stark überfordern. Die Voraussetzung zur Bewältigung der Vercodungsarbeit während des Interviews ist der Einsatz eines auf das äußerste reduzierten Codesystems, das - spontane - Befragtenverhalten nur noch im Hinblick darauf, ob es im Sinne der Fragestellung adäquat oder nicht adäquat ist, mittels einer Codeziffer im Fragebogen bewertet. Dabei bieten sich die beiden Ziffern „0“ für „adäquates Verhalten“ und „1“ für „nicht adäquates Verhalten“ an.

Eine weiteres Kennzeichen des Problem Coding liegt darin, daß der Interviewer im Unterschied zum traditionellen Behaviour Coding in einem zweiten Schritt zusätzlich bei inadäquater Verhaltensweise in einem schriftlichen Erfahrungsbericht nach dem Interview möglichst detailliert dieses Befragtenverhalten beschreibt. Damit erhält der Forscher Hinweise auf mögliche Ursachen für Mängel bei einer Frage.

Bei ZUMA hat sich die Problem Coding Technik in mehreren Studien bewährt. Für den Interviewer bedeutet der Einsatz der Technik allerdings eine hohe Anforderung, die nur durch entsprechende Schulungsmaßnahmen erfüllt werden kann.

2.1.3 Random Probe

Die "Random-Probe-Technik" wurde von Schuman (1966) mit dem Ziel entwickelt, das Fragenverständnis bei geschlossenen Fragen in Hauptstudien zu überprüfen. Dabei wählt jeder Interviewer vor dem Interview nach einem Zufallsverfahren eine bestimmte Anzahl von Fragen aus, bei denen Zusatzfragen (Probes) zum Fragenverständnis gestellt werden müssen. Beispielsweise wurden in der von Schumann erwähnten Studie pro Fragebogen jeweils zehn der insgesamt 200 Items zufällig ausgewählt. Als Probes standen dabei folgende drei Formulierungen zur Verfügung (Originaltext S. 241):

1. „Would you give me an example of what you mean?“
2. „I see - why do you say that?“
3. „Could you tell me a little more about that?“

In seiner Studie demonstriert Schuman an zwei Fragen die Eignung seiner Random-Probe-Technik und kommt zu folgendem Ergebnis (S. 244): „The answers to these questions show excellent variation, intercorrelate well, are significantly related to a number of background variables, and are relevant to an important hypothesis. But the random probes suggest that the questions were reasonably well understood by less than half the

sample.“ Obwohl die Random-Probe-Technik von Schuman ursprünglich zum Einsatz in Hauptstudien vorgesehen war, ist sie auch in Testerhebungen zur Evaluation von Fragen sinnvoll anwendbar.

2.1.4 Intensive Interview (Belson)

Belson (1981, 1986) kann nachweisen, daß auch formal korrekten Antworten ein falsches, d.h. vom Fragenkonstrukteur nicht intendiertes Verständnis des Frageinhalts zugrunde liegen kann. Mit der üblichen Preteststrategie, aus den Reaktionen bzw. Antworten der Befragten Rückschlüsse auf das Fragenverständnis zu ziehen, sind solche Fälle, bei denen trotz formal korrekter Antwort ein falsches Fragenverständnis vorliegt, nicht zu erkennen. Belson empfiehlt dafür ein Verfahren, bei dem Befragte nach der Durchführung eines Standard-Pretest-Interviews zum Verständnis von drei bis vier bereits vorher festgelegter Fragen intensiv befragt werden. Belson nennt dieses zweite Interview "Intensive Interview", das mittels eines zweistufigen Vorgehens erhoben wird:

1. Im ersten Schritt liest der Interviewer¹⁾ die zu testende Frage sowie die aus dem Standard-Pretest-Interview bereits vorliegende Antwort noch einmal vor. Der Befragte²⁾ wird anschließend gebeten, eine Beschreibung darüber zu geben, wie die Antwort zustande kam, wobei der Interviewer wie bei einem Tiefeninterview extensiv nachfragen soll.
2. Im zweiten Schritt stellt der Interviewer eine oder mehrere fest vorgegebene Fragen, um festzustellen, wie bestimmte Fragenaspekte beim vorangegangenen Standard-Pretest-Interview verstanden wurden.

Variationen dieser Technik sind auch unter den Bezeichnungen "Respondent Debriefing" (vgl. z.B. DeMaio 1996), „Reinterview“ (Bailar 1986), „Double Interview“ (Gordon 1963), „Intensive Reinterview“ (Johnson/Woltman 1986) oder „Follow-Up Interview“ (Morton-Williams/Sykes 1984) bekannt.

2.1.5 Qualitative Interviews

Unstrukturierte Interviews, Tiefeninterviews und ähnliche „qualitative“ Interviewformen können sinnvoll in einer frühen Entwicklungsphase des Fragebogens eingesetzt werden. Diese Interviewformen besitzen einen eher explorativen und experimentellen Charakter, d.h. sie dienen vorwiegend dazu, Ideen, Hinweise und Informationen zur Fragenkonstruktion zu generieren. Der Interviewer ist dabei von den Zwängen eines standardisierten Interviews befreit, d.h. er kann bei Bedarf nachfragen, hinterfragen oder alternative Frageversionen anbieten.

2.1.6 Analyse der Antwortverteilungen

Über die Häufigkeitsverteilung von Antwortalternativen lassen sich - meist nur grobe - Rückschlüsse auf die Qualität einer Frage ziehen. Indikatoren für Fragenmängel sind dabei in der Regel

- nicht oder nur minimal besetzte Antwort-Kategorien,
- extreme Häufigkeitsverteilung über die Antwort-Kategorien,
- hohe Häufigkeitswerte bei sog. „Ausweichkategorien“, wie z.B. „weiß nicht“ (Befragter kann sich nicht entscheiden oder hat keine Informationen) oder „verweigert“ (Befragter möchte die Frage nicht beantworten).

Sinnvoll ist dieses Verfahren nur bei einer genügend großen Fallzahl.

2.1.7 Split-Ballot

Beim Split-Ballot-Verfahren werden zwei (oder mehr) Varianten einer Frage jeweils einer Teilgruppe der Befragtenstichprobe zur Beantwortung präsentiert. Unterschiede in den Antwortverteilungen werden dann auf die unterschiedlichen Fragevarianten zurückgeführt.

Unter dem Aspekt der Evaluation von Fragen hat das Split-Ballot-Verfahren zum Ziel, eine Entscheidung für diejenige Fragenvariante herbeizuführen, die letztendlich zum Einsatz kommen soll. Diese Entscheidung trifft der Forscher normalerweise auf der Grundlage der Häufigkeitsverteilungen bzw. auf Grund von statistischen Analysen. Unter dieser Voraussetzung sollte ein Feld-Pretest, bei dem ein Split-Ballot-Verfahren eingesetzt wird, einen Stichprobenumfang von mindestens 100 Interviews haben. Neben Analyse- und Verteilungsaspekten können aber auch Pretestbeobachtungen als Entscheidungsgrundlage für eine bestimmte Formulierungsvariante einer Frage dienen. Dabei kann es sich um Pretestinformationen unter Einsatz eines traditionellen Standard-Pretests handeln, aber auch um Beobachtungen aus anderen Verfahren, wie z.B. Nachfaßfragen zum Verständnis bestimmter Frageninhalte (Probingverfahren).

2.2 Kognitive Laborverfahren

Aus der interdisziplinären Zusammenarbeit von Kognitionspsychologen und Umfrageforschern, deren Beginn auf das Ende der siebziger Jahre datiert werden kann, ging eine Reihe von Methoden hervor, die zwar nicht unbedingt neu waren, mit denen jedoch Informationen über kognitive Prozesse während des Frage-Antwort-Prozesses gesammelt werden können. Da diese Informationen Hinweise darüber geben, wie Befragte eine Frage bzw. bestimmte Elemente davon verstehen und interpretieren, sind sie damit auch zur Evaluation von Survey-Fragen geeignet. Diese Methoden sind in den letzten Jahren

unter der Bezeichnung „kognitive Laborverfahren“ bekannt geworden. Dabei muß das „Labor“ nicht unbedingt mit technischer Ausrüstung, wie z.B. Tonband, Videorecorder oder Einwegscheibe bestückt sein; in den meisten Fällen genügt ein schlichter Büroraum. In der Regel wird bei diesen Verfahren mit nur wenigen Befragten gearbeitet. Im folgenden sollen die wichtigsten kognitiven Laborverfahren kurz vorgestellt werden.

Think-Aloud

Diese Technik kann als die zentrale kognitive Technik überhaupt bezeichnet werden. Der Befragte wird aufgefordert, „laut zu denken“ und dabei sämtliche Gedankengänge, die zur Antwort führen bzw. führten zu formulieren. Ziel dabei ist, aus den Äußerungen der Befragten Hinweise darüber zu erhalten, wie die ganze Frage oder einzelne Begriffe verstanden wurden. Die - üblicherweise auf Tonträger - aufgezeichneten Formulierungen werden auch als "verbal protocols" (vgl. z.B. Ericsson/Simon 1980) bezeichnet. Über die Anwendung der Think-Aloud-Technik in der Umfrageforschung finden sich in der Literatur wenig klare Instruktionen. So weisen Blair/Presser (1993) anhand einer Befragung von 68 akademischen Institutionen in den USA nach, daß es beim Einsatz der Methode keine klaren Empfehlungen bezüglich Auswahl und Schulung der Interviewer, Anzahl der durchzuführenden Interviews, Bandaufzeichnungen und Analyseverfahren gibt. Bei der Anwendung der Think-Aloud-Methode gibt es zwei unterschiedliche Vorgehensweisen:

1. Die Befragten werden aufgefordert, laut zu denken, während sie ihre Antwort formulieren. Diese Vorgehensweise bezeichnet man als Concurrent-Think-Aloud-Methode.
2. Die Befragten werden aufgefordert, nach der Beantwortung der Frage zu beschreiben, wie die Antwort zustande kam. Diese Vorgehensweise ist bekannt unter dem Begriff Retrospektive-Think-Aloud-Methode.

Die Think-Aloud-Methode wurde zur Evaluation von Fragen in unterschiedlichen Bereichen erfolgreich eingesetzt:

1. Bei retrospektiven Fragen: Loftus (1984) setzte beispielsweise die Concurrent-Variante ein, um zu klären, wie Befragte bei der Frage, wie häufig sie in den letzten zwölf Monaten bei einem Arzt gewesen sind, vorgehen: Überlegen die Befragten vom gegenwärtigen Zeitpunkt ausgehend rückwärts oder umgekehrt, vom Zeitpunkt von vor zwölf Monaten bis in die Gegenwart? Die Concurrent-Think-Aloud-Methode konnte zeigen, daß bei autobiographischen Gedächtnisfragen Befragte eher in der „Vergangenheit-Gegenwart-Richtung“ denken. Ergebnisse eines Einsatzes der Concurrent-Think-Aloud-Methode bei ZUMA zeigten, daß bei retrospektiven Faktfragen zu Alltagsgeschehnissen, wie z.B. Fernsehkonsum der letzten sieben Tage, die Zeiten

weder vorwärts noch rückwärts aufaddiert werden, sondern in den meisten Fällen eine Schätzung des durchschnittlichen Verhaltens pro Tag zugrunde gelegt wird, um dieses dann für den entsprechenden Zeitraum hochzurechnen. Der Einsatz der Think-Aloud-Methode erweist sich demnach dann als sinnvoll, wenn es darum geht, den Antwortprozeß bei retrospektiven Faktfragen transparent zu machen. Kenntnisse dieser Art ermöglichen dann bessere und präzisere Formulierungen dieses Fragentyps.

2. **Bei Meinungsfragen:** In einer Studie der ZUMA-Feldabteilung wurde der Einsatz der Concurrent-Think-Aloud-Methode bei Meinungsfragen an 31 Fällen überprüft, wobei die Methode nicht wie üblich im Labor, sondern im Feld mit speziell geschulten Pretest-interviewern eingesetzt wurde. Dabei wurde unter anderem das in einer ALLBUS-Studie erhobene Item „Ein Mann schlägt sein 10-jähriges Kind, weil es ungehorsam war“ in die Überprüfung einbezogen. Das Item ist mittels einer 4-Punkte-Skala (sehr schlimm/ziemlich schlimm/weniger schlimm/überhaupt nicht schlimm) zu bewerten. Durch die Methode des lauten Denkens wurden Probleme der Befragten bei der Bewertung des Items deutlich. Es handelte sich dabei um die gleichen Probleme, die bereits bei der Durchführung eines Standard-Pretests bei diesem Item bekannt waren, nämlich eine zu starke Generalisierung der zu bewertenden Situation. Beim Einsatz der Concurrent-Think-Aloud-Methode traten die Probleme allerdings weit häufiger auf (in 29 Prozent aller Fälle, in denen die Befragten laut denken konnten), als dies beim Standard-Pretest der Fall war (2 Prozent).

3. **Zur Überprüfung von Hypothesen:** Bishop (1992) konnte nachweisen, daß sowohl die Concurrent- als auch die Retrospektive-Think-Aloud-Methode auch zur Hypothesenüberprüfung sinnvoll eingesetzt werden kann. Er wendet die Methode bei bereits bekannten Experimenten bezüglich Fragenabfolge und Kontexteffekten, wie z.B. dem bekannten Experiment von Schuman/Presser (1981) zu „Communist and American Reporters“ an, und weist nach, daß das, was Befragte bei ihrer Antwort laut dachten, genau dem entsprach, was Schuman/Presser als Erklärung des Kontexteffekts formulierten. Es gibt Befürworter für die eine und für die andere Vorgehensweise. So sprechen sich z.B. Sudman/ Bradburn/Schwarz (1996) für die retrospektive Variante aus, da Befragte erfahrungsgemäß den Prozeß, der zur Antwort führte, nicht immer in Worte fassen können. Die befragte Person wird dann im nachhinein, d. h. nach der Formulierung ihrer Antwort gebeten, ihre Überlegungen zu beschreiben, die zur Antwort führten.

Ähnliche Erfahrungen wurden auch in der ZUMA-Feldabteilung beim Einsatz der Concurrent-Think-Aloud-Methode gemacht. So waren bei zu skalierenden Meinungsfragen nur etwa die Hälfte der Befragten in der Lage, vor Nennung eines Skalenwertes ihre Ge-

danken laut zu formulieren, die letztendlich zu der Entscheidung für einen Skalenwert führten. Die Concurrent-Think-Aloud-Methode stellt an die Befragten hohe Anforderungen, die nur unter detaillierter Anleitung überhaupt erfüllt werden können.

Probing

Beim Probing handelt es sich um eine altbekannte Interview-Technik, die z.B. als zentrales Element Bestandteil der bereits beschriebenen Verfahren „Random Probe“ von Schuman (1966) und „Intensive Interview“ von Belson (1981) sind. Dabei wird eine gegebene Antwort vom Interviewer durch eine oder mehrere Zusatzfragen (Probes) „hinterfragt“, um mehr Informationen zu erhalten. Je nachdem, ob das Probing während des Interviews oder danach durchgeführt wird, werden folgende Bezeichnungen verwendet:

Follow-Up-Probing: Probing sofort nach der spontanen Antwort.

Post-Interview-Probing: Probing nach dem Interview.

Unabhängig vom Probing-Zeitpunkt kann auch nach der Aufgabenstellung, auf die sich das Probing bezieht, unterschieden werden. Hier werden z.B. von Oksenberg/Cannell/Kalton (1991) zwei weitere Probing-Varianten erwähnt:

Comprehension Probing: Probing zum Fragenverständnis.

Oksenberg/Cannell/Kalton (1991) nennen drei Varianten des Comprehension Probing:

1. Befragte sollen die Bedeutung eines bestimmten Begriffs in einer Frage erläutern.
2. Befragte sollen Aspekte ihrer Antwort erläutern.
3. Befragte sollen erläutern, wie klar verständlich ein Begriff für sie war oder welche Probleme sie beim Verständnis eines Begriffs hatten.

Information Retrieval Probing: Probing zu Aspekten der Informationsbeschaffung. Sinnvolles Anwendungsgebiet sind besonders retrospektive Faktfragen. Beispiel:

Frage: "Wann waren Sie zum letzten Mal beim Zahnarzt?" Information Retrieval Probing: "Wie schwer fiel es Ihnen, die Frage zu beantworten?"

Confidence Rating

Befragte sollen nach der eigentlichen Antwort den Grad der Verlässlichkeit ihrer Antwort bewerten (meist mit Hilfe einer Skala). Beispiel:

Frage: "Wie lange haben Sie in den letzten sieben Tagen insgesamt ferngesehen?" Confidence Rating: "Was würden Sie sagen: Ist Ihre Angabe sehr genau, ziemlich genau, eher ungenau oder grob geschätzt?"

Paraphrasing

Befragte sollen - nach der Beantwortung - die Frage mit eigenen Worten wiederholen bzw. formulieren. Erfahrungsgemäß gehen Befragte dabei unterschiedlich vor: Die einen versuchen, sich möglichst wörtlich an den Fragetext zu erinnern, die anderen versuchen, den Inhalt der Frage in eigenen Worten wiederzugeben. Drei Versuche von Befragten, den Text einer Frage zu wiederholen, sollen im folgenden als Beispiel angeführt werden.

Fragetext: "Im Vergleich dazu, wie andere hier in Deutschland leben: glauben Sie, daß Sie Ihren gerechten Anteil erhalten, mehr als Ihren gerechten Anteil, etwas weniger oder sehr viel weniger?" (Dabei wird die Skala optisch nicht vorgegeben.)

Nach Beantwortung der Frage werden die Befragten aufgefordert: "Bitte wiederholen Sie die Frage, die ich Ihnen eben vorgelesen habe noch einmal in Ihren eigenen Worten. Wie lautete die Frage?" Drei Antwortbeispiele aus kognitiven Laborinterviews, die in der ZUMA-Feldabteilung erhoben wurden und bei denen das Paraphrasing-Verfahren zum Einsatz kam, sollen die Wirkungsweise dieser Technik demonstrieren.

Formulierungsversuch eines Befragten: „Glauben Sie, daß Sie in Ihrer jetzigen Tätigkeit - verglichen mit anderen in Deutschland Lebenden - den gerechten Anteil bekommen, weniger gerecht, einigermaßen gerecht oder ganz ungerecht“.

Formulierungsversuch eines zweiten Befragten: Daß ich sagen sollte, daß ich im Vergleich zu anderen Bevölkerungsteilen über Maßen vom Sozialstaat profitiere“.

Formulierungsversuch eines dritten Befragten: „Ob ich eigentlich mit dem, was ich besitze, was ich habe, mit dem, was ich tun kann, zufrieden bin“.

Die Technik kann einem Forscher einerseits aufschlußreiche Hinweise geben, welche inhaltlichen Aspekte Befragte mit einer Frage verbinden, und andererseits kann die Paraphrasing-Technik zeigen, ob der Fragetext in allen Aspekten erinnert werden kann (z.B. konnte die 4-stufige Skala nicht korrekt wiedergegeben werden).

Sorting - Verfahren

Sorting - Verfahren sollen vornehmlich Hinweise darüber geben, wie Befragte Begriffe kategorisieren bzw. als Konzept verstehen. Es gibt drei Varianten:

1. Free Sort: Befragte sollen vorgegebene Items nach eigenen Kriterien gruppieren. Die Items werden dabei auf Kärtchen vorgegeben und sollen in selbstdefinierte Gruppen bzw. „Häufchen“ sortiert werden.

2. Dimensional Sort: Beim Dimensional Sort wird vorgegangen wie beim Free Sort, nur daß hier vorgegebene Items nach vorher festgelegten Kriterien sortiert werden sollen.

3. Vignette Classifications: Bei Vignette Classifications handelt es sich um eine Variante des Dimensional Sort. Beispielsweise sollen Befragte kurze Situationsbeschreibungen („Vignettes“) lesen und jeweils entscheiden, ob diese ihrer Meinung nach in die Überlegungen bei der Beantwortung einer vorgelegten Frage mit einbezogen werden sollen oder nicht.

Response Latency

Bei dieser Technik handelt es sich um die Messung der Zeit zwischen Präsentation der Frage und der Antwort. Die Möglichkeiten reichen dabei von exakter Messung (z.B. mittels Stoppuhr oder Computer) bis hin zu einer groben Schätzung durch den/die Testleiter/in mittels Kategorien, wie z.B. „kurz“, „mittel“, „lang“. Diese subjektiven Schätzungen werden auch als „Qualitative Timing“ bezeichnet. Lange „Reaktionszeiten“ werden dabei in der Regel als Indikator für Fragenmängel interpretiert.

2.2.1 Bewertung der kognitiven Laborverfahren

Als Gesamtbewertung aller hier dargestellten kognitiven Laborverfahren lassen sich folgende Vor- und Nachteile nennen.

Vorteile dieser Techniken

- Schnelle Durchführung
- Niedrige Kosten
- Die Techniken können innerhalb verschiedener Stadien der Fragebogenkonstruktion angewandt werden (z.B. kann sofort im Anschluß an eine Änderung einer Frage diese neue Version im Labor getestet werden).

Nachteile dieser Techniken

- Diese Techniken beschränken sich vorwiegend auf die Evaluation einzelner Fragen und nicht auf den Fragebogen als Ganzes. Dies bedeutet, daß diese Labor-Techniken keinesfalls einen Test des gesamten Instruments - in welcher Form auch immer - ersetzen können.
- Durch die geringe Fallzahl besteht ein hohes Unsicherheitsrisiko bezüglich der Generalisierbarkeit der Ergebnisse.

2.3 Andere Verfahren

Im Folgenden werden Evaluationsverfahren kurz vorgestellt, die weder der Kategorie „Testerhebung im Feld“ noch der Kategorie „kognitive Laborverfahren“ zugeordnet werden können.

2.3.1 Focus Groups

Im Bereich der Evaluation von Fragen können Focus Groups auf zwei Arten sinnvoll eingesetzt werden:

1. In einer frühen Entwicklungsphase des Fragebogens können Focus Groups wertvolle Hinweise zu Akzeptanz oder Verständnis des Befragungsthemas, einzelner Themenbereiche, einzelner Fragen oder einzelner Begriffe geben.
2. Focus-Groups eignen sich besonders dafür, schriftliche Fragebogen zu testen. Dabei bearbeitet zunächst jeder der Gruppenmitglieder für sich den Fragebogen, wobei keine Möglichkeit für Rückfragen gegeben werden sollte, da eine möglichst realistische Bearbeitungssituation simuliert werden soll, vor allem, um für jedes Gruppenmitglied die individuelle Bearbeitungsdauer festhalten zu können. Anschließend können die Gruppenmitglieder allgemeine Eindrücke zum Fragebogen, wie z.B. zur Thematik, zur Bearbeitungsdauer oder zum Schwierigkeitsgrad äußern. Danach wird der Fragebogen Frage für Frage „durchgearbeitet“, wobei die Gruppenmitglieder aufgefordert werden, zu jeder einzelnen Frage - soweit vorhanden - Kommentare, Verständnisprobleme oder Rückfragen zu äußern. Daneben werden vom Moderator an die Gruppe bereits vorbereitete Fragen zu einzelnen Fragen gestellt, überwiegend zum Verständnis der ganzen Frage oder einzelner Begriffe.

Grundsätzlich empfiehlt es sich, eine Focus-Group-Sitzung auf Tonträger aufzuzeichnen, als „Notlösung“ ist jedoch auch eine schriftliche Protokollführung durch einen Co-Moderator denkbar. Der Vorteil von Focus Groups liegt vor allem darin, daß mehrere Personen gleichzeitig „befragt“ werden können, ein entscheidender Nachteil ist darin zu sehen, daß soziale Interaktionen bzw. gruppenspezifische Prozesse, die zwangsläufig bei einer Focus-Group-Sitzung auftreten, das bei der eigentlichen Befragung relevante Individualverhalten verzerrt darstellen bzw. nicht adäquat repräsentieren. Diese Nachteile sowie die - meist - geringe Fallzahl lassen es ratsam erscheinen, den Fragebogen vor seinem endgültigen Einsatz einem Feld-Pretest unter möglichst realistischen Bedingungen zu unterziehen.

2.3.2 Experten

Zur Beurteilung von Fragebogen eines beliebigen Entwicklungsstadiums können Experten zu Rate gezogen werden. Dabei sollten idealerweise mehrere Experten eingesetzt werden, die ihre Bewertungen zur besseren Vergleichbarkeit anhand eines vorgegebenen Kriterienkatalogs vornehmen. Von Lessler und Forsyth (1996) wurde beispielsweise ein detailliertes Codesystem entwickelt, mit dessen Hilfe Experten eine Frage nach ihren

Merkmale und Eigenschaften - auch im Hinblick auf die Aufgabenstellung für die Befragten - beurteilen können („Expert Questionnaire Appraisal Coding System“).

3. Zusammenfassung und Ausblick

Bis Mitte der achtziger Jahre stand der Pretest nur äußerst selten im Blickpunkt des wissenschaftlichen Interesses. Er galt zwar in der älteren methodischen Literatur als wesentlicher Bestandteil im Gesamtkonzept einer Umfrage, gleichzeitig gab es wenig „übereinstimmende“ Anhaltspunkte in der Literatur für die konkrete Durchführung. In der Regel kam der Standard-Pretest zur Anwendung, obwohl man sich dessen Schwächen bewußt war und obwohl bereits 1966 Schuman zwecks besserer Information über das Verständnis von Fragen den Einsatz einer „Random Probe“ empfahl und Belson (1981, 1986) auf Grund seiner Studien die Notwendigkeit sah, formal adäquate Antworten der Befragten zu hinterfragen.

Heute stehen zur Evaluation von Fragen eine ganze Reihe von Verfahren zur Verfügung. Die sozialwissenschaftliche Methodenforschung, die im Bereich der Fragebogenkonstruktion durch die Zusammenarbeit mit Kognitionsforschern in den letzten Jahren zu äußerst praxisrelevanten Erkenntnissen kam, bezog ab Mitte der achtziger Jahre auch den Pretestbereich mit ein. Vor allem die amerikanische Literatur beschrieb den erfolgreichen Einsatz von „neuen“ Verfahren, wie z.B. „Think-Aloud“, „Probing“, oder „Paraphrasing“, die bislang entweder in anderen Forschungsbereichen zur Anwendung kamen oder einfach „in Vergessenheit“ geraten waren. Sie waren also nicht unbedingt neu, wurden aber wieder populär und für den Pretestbereich übernommen. Es sind die am häufigsten eingesetzten sogenannten kognitiven Laborverfahren.

Diese neuen kognitionspsychologischen Verfahren bieten den Vorteil, Einblick in die Gedankenprozesse der Befragten zu gewinnen, um so Probleme bei Fragen zu identifizieren. Im Gegensatz dazu ist die Identifizierung von Problemen beim Standard-Pretest ja nur dann der Fall, wenn Befragte selbst um Klärung bitten oder sich offensichtlich falsch verhalten.

Insbesondere hat der Einsatz solcher Verfahren dazu beigetragen, Erkenntnisse bei der Beantwortung retrospektiver Fragen zu gewinnen (vgl. z.B. Tanur 1992; Schwarz/Sudman 1994).

Die Vielzahl aktueller Evaluationsverfahren wirft die Frage auf, welche Verfahren sinnvoll eingesetzt werden können bzw. welche Verfahren für das jeweilige konkrete Umfrageprojekt geeignet sind. Neben den bereits genannten (und bekannten) Vor- und Nachteilen einzelner Verfahren geben insbesondere die Forschungsarbeiten von

Oksenberg/ Cannell/Kalton (1991) und Presser/Blair (1994), bei denen verschiedene Verfahren verglichen wurden, interessante Hinweise auf deren Leistungsfähigkeit.

Übereinstimmendes Fazit beider Arbeiten: Es gibt keine Methode, die in allen Problem-bereichen zufriedenstellend arbeitet. Oksenberg/Cannell/Kalton (1991) stellen in Ihrer Vergleichsstudie fest, daß sog. „Special Probes“ (z.B. Comprehension Probes) zwar erfolgreich zur Aufdeckung von Verständnisproblemen eingesetzt werden können, weniger jedoch zur Identifizierung aller anderen Probleme. Bewährt hat sich in dieser Studie auch das Behaviour Coding, wobei sich allerdings auch hier zeigte, daß die Ursachen der Probleme nicht direkt erkennbar sind. Presser/Blair (1994) berichten z.B., daß der Standard-Pretest im Vergleich zu anderen Verfahren am wenigsten reliabel ist. Im Gegensatz dazu ist das Behaviour Coding sehr reliabel auf Grund der Anwendung objektiver Regeln, es liefert aber keine Hinweise auf die Ursachen dieser Probleme. Kognitive Verfahren wie Probes und Think-Aloud-Verfahren liefern die meisten Verständnisprobleme, aber z.B. keine Interviewerprobleme. Das Verfahren der Expertenrunde liefert vergleichsweise die meisten Erkenntnisse und ist am kostengünstigsten, besitzt allerdings starke Defizite bei Hinweisen auf Interviewerprobleme.

Auf der Grundlage dieser Ergebnisse empfiehlt es sich also, mehrere Verfahren einzusetzen. Da der Erkenntniswert der einzelnen Verfahren für unterschiedliche Problembe-reiche differiert, sollte der Einsatz der Verfahren sinnvoll kombiniert werden. So emp-fiehlt beispielsweise Fowler (1995) für die Evaluation von Fragen den Einsatz von Focus Groups, kognitiven Laborinterviews und einen Feld-Pretest mit Auswertung der Ant-wortverteilungen. Unabhängig davon, welche Verfahren kombiniert eingesetzt werden, sollte auf einen abschließenden Feld-Pretest auf keinen Fall verzichtet werden, da nur hier wichtige Informationen über Interviewerprobleme oder die Wirkungsweise des ge-samten Fragebogens (wie z.B. die Befragungsdauer, Sukzessionseffekte) gesammelt wer-den können.

Anmerkungen

- 1) Der Einfachheit halber wird im Text immer nur die männliche Form verwendet.
- 2) Der Einfachheit halber wird im Text immer nur die männliche Form verwendet.

Literatur

Attleslander, P., 1984: Methoden der empirischen Sozialforschung. Berlin: Walter de Gruyter.

- Bailar, B. A., 1986: Recent Research in Reinterview Procedures. S. 41 - 63 in: *Journal of the American Statistical Association*, 63, 1986.
- Belson, W. A., 1981: *The Design and Understanding of Survey Questions*. Aldershot, England: Gower.
- Belson, W. A., 1986: *Validity in Survey Research*. Aldershot, England: Gower.
- Bishop, G., 1992: Qualitative Analysis of Question-Order and Context Effects: The Use of Think-Aloud Responses. S. 149-162 in: N. Schwarz/S. Sudman (Hrsg.), *Context Effects in Social and Psychological Research*. New York: Springer.
- Blair, J./Presser, S., 1993: *Survey Procedures for Conducting Cognitive Interviews to Pretest Questionnaires: A Review of Theory and Practice*. Survey Research Center, University of Maryland.
- Bolton, R., 1993: Pretesting Questionnaires: Content Analyses of Respondents' Concurrent Verbal Protocols. S. 280-303 in: *Marketing Science*, 12, 1993.
- Cannell, C. F./Axelrod, M., 1956: The Respondent Reports on the Interview. *American Journal of Sociology*, 62, 1956.
- Cannell, C./Lawson, S./Hausser, D., 1975: *A Technique for Evaluating Interviewer Performance*. Ann Arbor: The University of Michigan; Survey Research Center; Institute for Social Research.
- Cannell, C./Kaltan, G./Fowler, F., 1985: *Techniques for Diagnosing Cognitive and Affective Problems in Survey Questions*. Ann Arbor: The University of Michigan; Survey Research Center; Institute for Social Research.
- Cannell, C./Oksenberg, L./Kaltan, G./Bischoping, K./Fowler, F. J., 1989: *New Techniques for Pretesting Survey Questions. Final Report*. Ann Arbor: The University of Michigan; Survey Research Center. Boston: University of Massachusetts; Center for Survey Research.
- Converse, J. M./Presser, S., 1986: *Survey Questions. Handcrafting the Standardized Questionnaire*. Beverly Hills: Sage.
- DeMaio, Th.(Hrsg.), 1983: *Approaches to Developing Questionnaires*. Bureau of the Census. Office of Management and Budget. Statistical Working Paper 10.
- DeMaio, Th./Rothgeb, J. M., 1996: *Cognitive Interviewing Techniques: In the Lab and in the Field*. S. 177-195 in N. Schwarz/S. Sudman (Hrsg.), *Answering Questions*. San Francisco: Jossey-Bass.
- Elliott, K./Christopher, M., 1973: *Research Methods in Marketing*. London: Holt, Rinehart & Winston.
- Ericsson, K. A./Simon, H. A., 1980: Verbal Reports as Data. S. 215 - 251 in: *Psychological Review*, 8, 1980.
- Forsyth, B. H./Lessler, J. T., 1991: *Cognitive Laboratory Methods: A Taxonomy*. S. 393 - 418 in: P. P. Biemer/R. M. Groves/L. E. Lyberg./N. Mathiowetz/S. Sudman (Hrsg.): *Measurement Errors in Surveys*. New York: Wiley.

- Fowler, F. J., 1984: *Survey Research Methods*. Beverly Hills: Sage.
- Fowler, F. J., 1992: How unclear Terms Affect Survey Data. S. 218 - 231 in: *Public Opinion Quarterly*, 56, 1992.
- Fowler, F. J., 1995: *Improving Survey Questions*. Thousand Oaks: Sage.
- Friedrichs, J., 1973: *Methoden empirischer Sozialforschung*. Reinbek: Rowohlt.
- Gordon, W. D., 1963: Double Interview. In: *New Developments in Research*. London: Market Research Society with the Oakwood Press.
- Hunt, S. D./Sparkman, R. D./Wilcox, J. B., 1982: The Pretest in Survey Research: Issues and Preliminary Findings. S. 269 - 273 in: *Journal of Marketing Research*, Vol. XIX, 1982.
- Jabine, T. B./Straf, M. L./Tanur, J. M./Tourangeau, R. (Hrsg.), 1984: *Cognitive Aspects of Survey Methodology: Building a Bridge Between Disciplines*. Washington, D. C.: National Academy Press.
- Jobe, J. B./Mingay, D. J., 1989: Cognitive Research Improves Questionnaires. S. 1053 - 1055, in: *American Journal of Public Health*, 79, 1989.
- Jobe, J. B./Mingay, D. J., 1990: Cognitive Laboratory Approach Designing Questionnaires for Surveys of the Elderly. S. 518 - 524 in: *Public Health Reports*, 105, 1990.
- Johnson, R. A./Woltman, H. F., 1986: Evaluating Census Data Quality Using Intensive Reinterviews: A Comparison of U.S. Census Methods and Rash Methods. S. 293 - 298 in: *Proceedings of the Section on Survey Research*, American Statistical Association, 1986.
- Karmasin, F./Karmasin, H., 1977: *Einführung in die Methode und Probleme der Umfrageforschung*. Wien: Hermann Böhlau Nachf.
- Kidder, L. H., 1981: *Research Methods in Social Relations*. New York: Holt, Rinehart and Winston.
- Kreiselmaier, J./Prüfer, P./Rexroth, M., 1989: *Der Interviewer im Pretest*. Mannheim: ZUMA-Arbeitsbericht 89/14.
- Krueger, R. A., 1988: *Focus Groups. A Practical Guide for Applied Research*. Newbury Park: Sage.
- Lessler, J. T./Forsyth, B. H., 1996: A Coding System for Appraising Questionnaires. S. 259-291 in: N. Schwarz/S. Sudman (Hrsg.), *Answering Questions*. San Francisco: Jossey-Bass.
- Loftus, E., 1984: Protocol Analysis of Responses to Survey Recall Questions. In: T. B. Jabine/M. L. Straf/J. M. Tanur/R. Tourangeau (Hrsg.), 1984: *Cognitive Aspects of Survey Methodology: Building a Bridge Between Disciplines*. Washington, D. C.: National Academy Press.
- Morgan, D. (Hrsg.), 1993: *Successful Focus Groups*. Newbury Parks: Sage.
- Morton-Williams, J., 1979: The Use of "Verbal Interaction Coding" for Evaluating a Questionnaire. *Quality and Quantity*, 13, 1979: 59 - 75.

-
- Morton-Williams, J./Sykes, W., 1984: The Use of Interaction Coding and Follow-up Interviews to investigate Comprehension of Survey Questions. S. 109 - 127 in: Journal of the Market Research Society, 26, 1984.
- Noelle, E., 1971: Umfragen in der Massengesellschaft. Reinbek: Rowohlt.
- Oksenberg, L./Cannell, Ch./Kalton, G., 1991: New Strategies for Pretesting Survey Questions. Journal of Official Statistics, 7: 349 - 365.
- Oppenheim, A. N., 1966: Questionnaire Design and Attitude Measurement. New York: Basic Books.
- Payne, S. L., 1951: The Art of Asking Questions. Princeton, N.J.: Princeton University Press.
- Porst, R., 1985: Praxis der Umfrageforschung. Stuttgart: Teubner.
- Presser, S./Blair, J., 1994: Survey Pretesting: Do different Methods produce different Results? Sociological Methodology: 73 - 104.
- Prüfer, P./Rexroth, M., 1985: Zur Anwendung der Interaction-Coding-Technik. ZUMA-Nachrichten, 17: 2-49.
- Prüfer, P./Rexroth, M., 1996: Multi-Method-Pretesting. Manuskript. ZUMA, Mannheim.
- Reynolds, N./Diamantopoulos, A./Schlegelmilch, B., 1993: Pretesting in Questionnaire Design: A Review of the Literature and Suggestions for Further Research. Journal of the Market Research Society, 35, Nr. 2: 171 - 182.
- Royston, P./Bercini, D./Sirken, M./Mingay, D., 1986: Questionnaire Design Research Laboratory. S. 703 - 706 in: Proceedings of the Section on Survey Research, American Statistical Association, 1986.
- Schnell, R./Hill, P. B./Esser, E., 1995: Methoden der empirischen Sozialforschung. München/Wien: Oldenbourg.
- Schrader, A., 1971: Einführung in die empirische Sozialforschung. Ein Leitfaden für die Planung, Durchführung und Bewertung von nicht-experimentellen Forschungsprojekten. Stuttgart: Kohlhammer.
- Schuman, H., 1966: The Random Probe: A Technique for Evaluating the Validity of Closed Questions. American Sociological Review, 31: 218 - 222.
- Schuman, H./Presser, S., 1981: Questions and Answers in Attitude Survey: Experiments on Question Form, Wording and Context. New York: Academic Press.
- Schwarz, N./Sudman, S. (Hrsg.), 1992: Context Effects in Social and Psychological Research. New York: Springer.
- Schwarz, N./Sudman, S. (Hrsg.), 1994: Autobiographical Memory and the Validity of Retrospective Reports. New York: Springer.
- Schwarz, N./Sudman, S. (Hrsg.), 1996: Answering Questions. San Francisco: Jossey-Bass.

- Sheatsley, P. B., 1983: Questionnaire Construction and Item Writing. In: Rossi, P. H./Wright, J. D./Anderson, A. B. (Hrsg.): Handbook of Survey Research. New York: Academic Press.
- Sudman, S./Bradburn, N., 1982: Asking Questions. A Practical Guide to Questionnaire Design. San Francisco: Jossey-Bass.
- Sudman, S./Bradburn, N./Schwarz, N., 1996: Thinking About Answers. The Application of Cognitive Processes to Survey Methodology. San Francisco: Jossey-Bass.
- Tanur, J. M. (Hrsg), 1992: Questions about Questions. New York: Russell Sage Foundation.
- Warwick, D. P./Lininger, C. A., 1975: The Sample Survey: Theory and Practice. New York: Mc Graw - Hill.
- Wellenreuther, 1982: Grundkurs: Empirische Forschungsmethoden: Königstein: Athenäum.
- Williamson, J./Karp, D./Dalphin, J.R., 1977: The Research Craft: An Introduction to Social Science Methods. Boston: Little, Brown and Co.
- Willis, G. B./Royston, P./Bercini, D., 1991: The Use of Verbal Report Methods in the Development and Testing of Survey Questionnaires. Applied Cognitive Psychology, 5: 251 - 267.