

# **Computerunterstützte Inhaltsanalyse: Codierung und Analyse von Antworten auf offene Fragen**

Cornelia Züll & Peter Ph. Mohler

Zentrum für Umfragen, Methoden und Analysen, Mannheim

## *Zusammenfassung*

*Das Vorgehen bei der Verwendung der computerunterstützten Inhaltsanalyse zur Codierung von Antworten auf offene Fragen eines Fragebogens wird dargestellt. Der Schwerpunkt liegt dabei auf der Beschreibung der Entwicklung und der Validierung eines Kategoriensystems/ Diktionärs zur automatischen Codierung.*

## *Summary*

*The paper discusses the use of computer-assisted content analysis for coding open-ended answers in a questionnaire. The emphasis is on development and validation of a categorization scheme/dictionary for automatic coding.*

**ZUMA How-to-Reihe, Nr. 8**

**2001**

## 1. Einleitung

Das Papier beschreibt die Nutzung der computerunterstützten Inhaltsanalyse zur Codierung und Analyse von Antworten auf offene Fragen, die im Rahmen einer Umfrage als Teil eines Fragebogens mit vorwiegend geschlossenen Fragen erhoben wurden. Anhand eines Beispiels wird die Vorgehensweise bei der automatischen Codierung dargestellt. Sie ist in weiten Teilen auch auf die Codierung anderer Textarten übertragbar. Obwohl es im Rahmen der computerunterstützten Inhaltsanalyse verschiedene Ansätze gibt, werden wir uns auf die Anwendung des klassischen diktionär-basierten Ansatzes beschränken.

Als Beispieldaten für die Codierung von Antworten auf offenen Fragen werden die Daten der Nachwahlstudie zur Bundestagswahl 1994<sup>1</sup> verwendet. An der Umfrage beteiligten sich 1000 Befragte aus West- und 1046 Befragte aus Ostdeutschland. Die Teilnehmer an der Umfrage wurden im Rahmen eines Fragebogens in geschlossenen Fragen zu ihrem politischen Engagement, nach der Partei, die sie gewählt haben, nach ihrer bevorzugten Partei, nach ihrer Einbindung in das politische System, nach ihrer Mediennutzung usw. befragt. Einige wenige Fragen des Fragebogens waren als offene Fragen formuliert. In den folgenden Beispielen werden wir die folgenden offenen Fragen verwenden:

- Wenn Sie an die ehemalige DDR zurückdenken, gibt es Dinge auf die die Menschen dort stolz sein können? Sie wurde an die westdeutschen Befragten gerichtet.
- Wenn Sie an die ehemalige DDR zurückdenken, gibt es Dinge auf die Sie stolz sind? Sie wurde an die ostdeutschen Befragten gerichtet.

## 2. Konventionelle vs. computerunterstützte Inhaltsanalyse

Als konventionelle oder intellektuelle Inhaltsanalyse werden Verfahren dann bezeichnet, wenn die jeweiligen Textmerkmale mit Hilfe „menschlicher“ Vercoder identifiziert werden, d.h. die Verfahren der konventionellen Inhaltsanalyse erfordern immer, daß Vercoder Teile des Textes bestimmten Merkmalen (Kategorien) zuordnen. Bei gut trainierten Coderteams können mit konventionellen inhaltsanalytischen Verfahren auch hochkomplexe Vercodungen zuverlässig durchgeführt werden. Allerdings steigen die Kosten mit der Komplexität der Kategorien und, unter Umständen, auch der Textmenge exponentiell an. Dies kann u.a. zur Folge haben, daß eine einmal begonnene Vercodung alleine aus ökonomischen Gründen fortgeführt wird, obwohl auf Grund der im Vercodungsverlauf gewonnenen Erkenntnisse von neuem begonnen werden müßte (Früh 1998).

In der computerunterstützten Inhaltsanalyse werden im Gegensatz zur konventionellen Inhaltsanalyse die Textmerkmale mit Hilfe von Computerprogrammen identifiziert. Unter automatischem Codieren versteht man also das Zuweisen von Codes, die in Diktionären definiert sind, zu Texteinheiten durch Computerprogramme. Nach dieser Definition ist eine Verarbeitung von manuell erzeugten Merkmalshäufigkeiten durch statistische Programme keine computerunterstützte Inhaltsanalyse. Eine solche automatische Codierung ist nur dann möglich, wenn die Vercodungsregeln explizit als logische Bedingungen (z.B. wenn-dann) formuliert und in Algorithmen, d.h. Computerbefehle, gefaßt werden. Die Umsetzung bestimmter theoretisch abgeleiteter Operationalisierungen macht die computerunterstützte Inhaltsanalyse spannend und schwierig zugleich. Spannend, indem komplexe theoretische Operationalisierungen Stück für Stück auf das logisch Mögliche und Notwendige reduziert werden; schwie-

---

<sup>1</sup> Die Studie "Nachwahlstudie zur Bundestagswahl 1994" ist unter der Studiennummer 2601 beim Zentralarchiv in Köln erhältlich.

rig, indem nicht immer sofort, manchmal auch gar nicht, die geforderte logische Klarheit und Eindeutigkeit gefunden werden können.

### 3. Codierung von Antworten auf offenen Fragen

#### 3.1 Die Texte

Jede Inhaltsanalyse setzt voraus, daß nach Bestimmung des Untersuchungsziels und der Hypothesen die Texte ausgewählt werden, die zur Analyse verwendet werden sollen.

Will man den Geltungsbereich der computerunterstützte Inhaltsanalyse bezüglich geeigneter Texte angeben, dann gilt, daß, bezogen auf die jeweilige Forschungsfrage, in den Texten genügend semantische Redundanz vorhanden sein sollte, um bei der weiteren Analyse quantifizieren, d.h. zählen, zu können. Weiterhin muß diese semantische Redundanz als Häufigkeit von Wortkategorien abbildbar sein. Anders gesagt, für die computerunterstützte Inhaltsanalyse muß man Textwörter so geschickt zu Kategorien zusammenfassen, daß aus der Anordnung oder dem Muster der Kategorien auf den Textinhalt geschlossen werden kann.

Festzulegen sind neben der Art der Texte (z.B. Zeitungsartikel, Anzeigen, e-mails), die Grundgesamtheit und die zu ziehende Stichprobe. Dieses Papier beschränkt sich auf die Inhaltsanalyse von Antworten auf offene Fragen, bei denen diese Auswahl durch die Integration in eine Umfrage bereits vorgegeben ist:

- Bei den Texten handelt es sich um Antworten auf offene Fragen, die in einer Umfrage erhoben wurden. Die offenen Fragen sind Teil eines ansonsten standardisierten Fragebogens.
- Die Grundgesamtheit und gegebenenfalls die Stichprobe werden durch die Umfrage selbst festgelegt.

Die Antworten auf die offenen Fragen (siehe oben) müssen zunächst verschriftet und als eigener Textkorpus, d.h. als Textdatei, gespeichert werden. Die Antworten sind häufig sehr kurz, bestehen oft nur aus einem einfachen Satz, einer Phrase oder einem oder mehreren Stichwörtern. Beispiele für solche Texte als Antwort auf die Frage nach den Objekten des Stolzes in der ehemaligen DDR finden sich in Tabelle 1.

Tabelle 1: Antworten der Befragten auf die offenen Fragen („\$“ kennzeichnet die Fragebogennummer)

\$197

Gutes Gesundheitssystem; Kindergärten kostenlos

\$198

Auf die Leistungen der Sportler

\$476

Sportleistungen

\$477

Kultur und Landschaft

\$478

Landschaft; Sportleistungen

\$479

Kultur

\$3234

Das Sozialwesen funktionierte gut. Die Arbeitsplatzsicherung war ganz wichtig.

\$3235

In der damaligen DDR gab es keine Arbeitslosigkeit. Das Sozialsystem war gut ausgebaut. Polizeigesetz war straffer und es wurde auch härter durchgegriffen.

\$3236

Die Kindererziehung (Kindergarten und Kinderkrippe) war sozial abgesichert. Die Arbeitsplätze waren gesichert. Wohnungsnot war nicht so schlimm wie jetzt.

§3237

Es gab keine Arbeitslosigkeit, die Zukunft war gesichert, keine Angst um den Arbeitsplatz. Die soziale Absicherung war da. In der ehemaligen DDR war das Sozialsystem gut (Kinderpflege), gab es für alle Kindergartenplätze und Hortplätze. Man macht die DDR manchmal zu schlecht.

§3238

In der ehemaligen DDR war das Sozialsystem gut (Krankenpflege), gab es für alle Kindergartenplätze. Man macht die DDR manchmal zu schlecht.

Bei jeder Inhaltsanalyse – konventionell oder computerunterstützt durchgeführt - besteht ein wichtiger Schritt darin, das Material zu sichten und die Texte zu lesen, um sich einen ersten Eindruck der Texte zu verschaffen. Beim Betrachten der hier vorliegenden Texte wird deutlich, daß sich die Texte der westdeutschen Befragten nicht nur durch die Zahl der vorliegenden Antworten (288 westdeutsche und 794 ostdeutsche Befragte haben die Frage beantwortet), sondern auch durch die Struktur der Antworten unterscheiden. Die westdeutschen Befragten (Tabelle 1, Antwort 1-5) antworteten meist mit einem Stichwort oder einer kurzen Phrase (im Durchschnitt mit 9,1 Wörtern pro Antwort), die ostdeutschen nannten dagegen oft eine ganze Liste von Bereichen (durchschnittlich 12,75 Wörter pro Antwort). Die Ausführlichkeit der Antworten ist ein erster Indikator für die Relevanz und Wichtigkeit des Themas für die Befragten.

Nach der Aufbereitung der Texte muß in einem weiteren Schritt ein geeignetes Diktionär aufgebaut oder ein bereits bestehendes Diktionär adaptiert werden.

### **3.2 Aufbau eines Diktionärs zur Codierung**

Als Programm für die computerunterstützte Inhaltsanalyse wird im folgenden das von ZUMA entwickelte Programm TEXTPACK (Mohler & Zuell 1995, <http://www.gesis.org/software/textpack/index.htm>) verwendet. Die Formate des Diktionärs und die Beispielausgaben sind jeweils TEXTPACK-spezifisch, das hier beschriebene Vorgehen ist jedoch programm-unabhängig.

Kernstück einer Inhaltsanalyse sind Kategorien, die bei der computerunterstützten Inhaltsanalyse in einem Diktionär (Wörterbuch) definiert werden. Der Aufbau eines Diktionärs kann sowohl theoriegeleitet wie auch empiriegeleitet erfolgen. In der Regel wird man sich für eine Kombination der beiden Verfahren entscheiden und nach einer zunächst theoriegeleiteten Herangehensweise das Diktionär empiriegeleitet erweitern.

### **3.3 Theoriegeleitete Kategorienbildung**

Das theoriegeleitete Vorgehen bei der Definition von Kategorien erfolgt zunächst völlig unabhängig von den vorliegenden Texten. Aus der Theorie und den Hypothesen lassen sich bereits Begriffe ableiten, die die künftigen Kategorien bilden sollen.

Hypothesen wie „Soziale Aspekte wie z.B. Arbeitsplatzsicherheit, Kindergartenplätze, soziale Absicherung sind für sowohl ostdeutsche wie auch westdeutsche Befragte das wichtigste Thema, auf das man in der ehemaligen DDR stolz sein konnte.“ oder „Politische Themen (Demokratie, Parteien, usw.) spielen für die Befragten keine Rolle.“ geben bereits einige Kategorien vor, die zu definieren sind:

- Kategorie „Arbeit,“ die Arbeitsplatzsicherheit, ausreichend Arbeitsplätze, Vollbeschäftigung einschließt;
- Kategorie „Soziale Dienste“, die Kindergarten- und Hortplätze genauso einschließt wie Altenheime;
- Kategorie „Demokratie“, die u.a. die Parteien und freie Wahlen beinhaltet.

Dies sind nur einige Beispiele, wie aus den Hypothesen erste Kategoriendefinitionen erfolgen können.

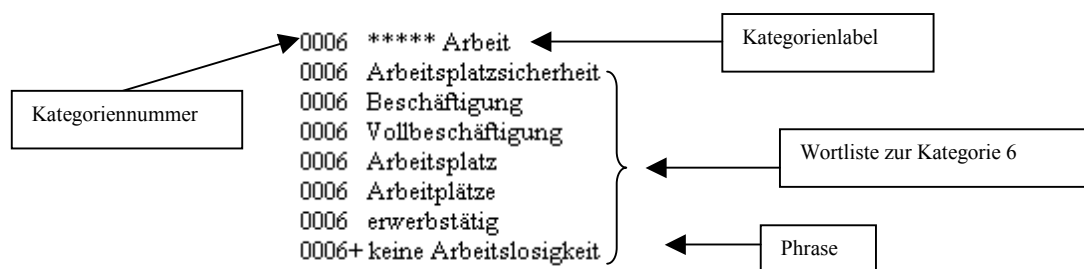
Kategoriendefinitionen unterliegen einigen Regeln:

- Kategorien müssen für die Forschungsfrage *relevant* sein, d.h. wenn man - wie in unserem Beispiel – die Frage nach der Meinung zu Leistungen der ehemaligen DDR analysieren will, ist es nicht sinnvoll Kategorien wie z.B. Medieninhalte zu definieren.
- Kategorien müssen *vollständig* definiert werden, d.h. alle Themen müssen berücksichtigt werden. Um Vollständigkeit zu erreichen, empfiehlt sich in der konventionellen Inhaltsanalyse die Definition der Restkategorie „Sonstiges“, in der computerunterstützten Inhaltsanalyse gehören alle nichtcodierbaren Nennungen zu dieser Restkategorie.
- Kategorien müssen *eindeutig* sein, d.h. sie müssen so definiert werden, daß z.B. Wörter nicht Indikatoren für mehrere Kategorien sind. Z.B. muß die Definition der Kategorie eindeutig festlegen, ob „Frauen“ ein Indikator für die Kategorie „Soziale Beziehungen“ oder für die Kategorie „Frauenrechte und Abtreibung“ ist (siehe Reliabilität und Validität).
- Kategorien müssen *trennscharf* sein, d.h. sie müssen voneinander unabhängig sein. Es ist z.B. nicht sinnvoll, neben einer Kategorie „Kultur und Bildung“ einen weitere Kategorie „Theater“ in das Diktionär aufzunehmen.

Kategorien werden bei computerunterstützten Inhaltsanalysen durch Wortlisten bzw. Listen von Mehrwortverbindungen (Phrasen) definiert. Eine Mehrwortverbindung setzt sich dabei aus mehreren Einzelwörtern zusammen (z.B. „keine Arbeitslosigkeit“, „soziale Sicherheit“) und ist in TEXTPACK mit einem „+“-Zeichen gekennzeichnet. In TEXTPACK wird jeder Kategorie ein numerischer Wert (ein Code) zugeordnet, der bei späteren statistischen Analysen verwendet werden kann. Die Kategorie „Arbeit“ erhält nun z.B. die Kategoriennummer „6“ und das Kategorienlabel „Arbeit“. Das Label wird in der späteren Analyse mit SPSS oder SAS als Variablenlabel verwendet. Alle Wörter, die die Kategorie definieren und in einer Liste zusammengestellt werden, sind eindeutige Indikatoren dafür, daß die Textstelle, in der sie auftreten, zu dieser Kategorie gehört. Der Aufbau eines Diktionärs erfordert also:

- begriffliche Definition der Kategorien;
- Erstellen der Wortliste und der Liste der Mehrwortverbindungen, durch die die Kategorien definiert sind;
- Vergabe von Kategorienlabels;
- Festlegen von Kategoriennummern (Codes).

Aus der oben beschriebenen Definition kann für die Kategorie „Arbeit“ der folgende Diktio-näreintrag (Wortliste) definiert werden:



### 3.4 Empiriegeleitete Kategorienbildung

Nachdem ein erster Grundstein für das Diktionär unabhängig vom Text gelegt wurde, können bei der empiriegeleiteten Vorgehensweise aus dem Text heraus weitere Kategorien definiert, bzw. die bereits bestehenden Kategorien erweitert werden. Worthäufigkeitslisten und Key-word-In-Context (KWIC)-Listen, die weiter unten erläutert werden, sind dazu wichtige Hilfsmittel.

Eine Liste aller Wörter eines Textes mit den dazugehörigen Häufigkeiten hilft, Themen aufzufinden, die für die Befragten wichtig sind. Sie hilft aber auch, Themen zu finden, die bei den Befragten überhaupt keine Rolle spielen. In unserem Beispiel ist ein solches Thema „Demokratie“, das nur dreimal genannt wird. Die Wortlisten können reduziert und damit übersichtlicher werden, indem man Wörter definiert, die für die Analyse keine Rolle spielen (Stop-Wörter) und die nicht angezeigt werden sollen (z.B. Artikel, Konjunktionen, Präpositionen). Die Listen können zudem verschieden sortiert werden, z.B. alphabetisch oder nach Worthäufigkeiten (die häufigsten Wörter am Anfang). Ein Auszug aus der Worthäufigkeitsliste unseres Textes ist in Tabelle 2 wiedergegeben.

Aus dieser Liste ergeben sich im vorliegenden Beispiel nun weitere Kategorien (z.B. „Sport“, „Kultur“), die in das Diktionär aufgenommen werden müssen, bzw. es werden Wörter aufgelistet, die zusätzlich zur Beschreibung einer bereits vorhandenen Kategorie in das Diktionär aufgenommen werden müssen (z.B. Betreuung als Indikator für die Kategorie „Soziale Dienste“).

Tabelle 2: Auszug aus einer Worthäufigkeitsliste

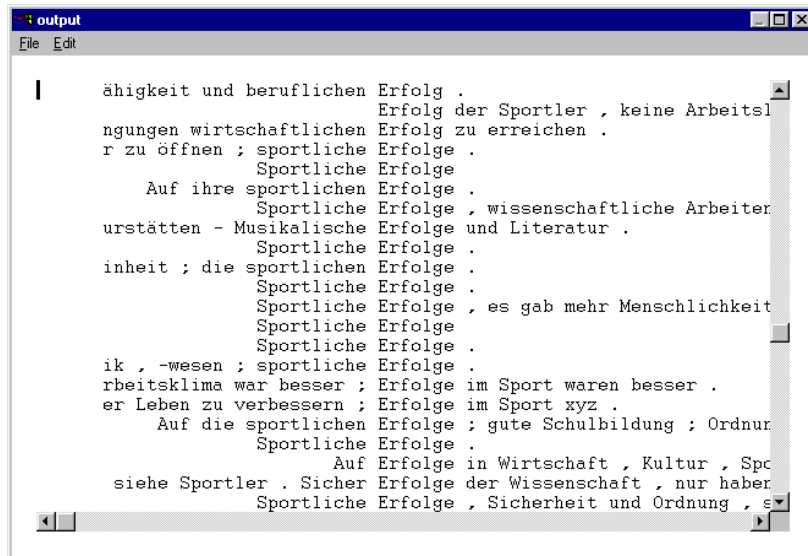


Wort	Häufigkeit
Kinderbetreuung	46
Sozialleistungen	44
Mieten	43
sozialen	42
Gesundheitswesen	41
Ordnung	39
Sozialwesen	36
Zusammenhalt	34
Kultur	33
Frauen	32
Sport	30
untereinander	30
Maßnahmen	29
Sozialsystem	29
Wohnung	29
Frau	28
Betreuung	26
Kindergartenplätze	26
System	26
Einrichtungen	23
Erfolge	23
Jeder	23

Bei einigen Wörtern wird man nicht sicher sein, ob sie einer der vorab definierten Kategorien zugeordnet werden können. In diesem Fall helfen KeyWord-In-Context-Listen (KWIC-Listen). Diese stellen Wörter im Satzzusammenhang dar. Sie helfen bei der Entscheidung, ob das Wort ein eindeutiger Indikator für eine bestimmte Kategorie ist oder ob ein Wort z.B. nur innerhalb einer bestimmten Mehrwortverbindung ein Indikator für eine der vorab definierten Kategorien ist. Eine KWIC-Liste zeigt das ausgewählte Wort in der Mitte einer Zeile, rechts und links des Wortes wird der Text (Kontext), der vor bzw. hinter dem Wort steht angezeigt. Der Kontext kann, wenn er für eine Feststellung des Gebrauch eines Wortes nicht ausreicht, auf die ganze Antwort erweitert werden. Ein Beispiel für eine solche KWIC-Liste ist in Tabelle 3 wiedergegeben und zeigt den Kontext der Wörter „Erfolg“ und „Erfolge“. In diesem

Auszug sieht man, daß die beiden Wörter allein kein eindeutiger Indikator für eine der vorher festgelegten Kategorien sind: z.B. muß „Musikalische Erfolge“ der Kategorie „Kultur“, „sportliche Erfolge“ dagegen der Kategorie „Sport“ zugeordnet werden. Anstelle der einzelnen Wörter werden die entsprechenden Mehrwortverbindungen in das Diktionär übernommen.

Tabelle 3: Auszug aus einer KWIC-Liste



Nachdem alle Kategorien definiert und das Diktionär erstellt wurde, sind die in Tabelle 4 dargestellten Kategorien im Diktionär zusammengefaßt.

Tabelle 4: Auszug aus dem Diktionär

Kategorie 0001 „Soziale Sicherheit allgemein“

Beispiele:

0001 Grundbedürfnisse

0001 Grundsicherung

0001 Absicherung

Zwei weitere Kategorien beschreiben spezielle soziale Sicherheit:

Kategorie 0011 „Soziale Dienste“

Beispiele:

0011 Kindergarten

0011 Hort

0011 Ferienplätze

Kategorie 0012 „Medizinische Versorgung“

Beispiele:

0012 Arzt

0012 Krank...

0012 Medikamente

Kategorie 0002 „Kosten und Preise“

Beispiele:

0002 Günstiges

0002 Strompreis

0002 preiswert

- Kategorie 0003 „Sport“  
 Beispiele:  
 0003 Spitzensportler  
 0003 Sport
- Kategorie 0004 „Soziale Beziehungen“  
 Beispiele:  
 0004 Familie  
 0004 Freunde  
 0004 Gemeinschaftssinn
- Kategorie 0005 „Frauenrechte und Abtreibung“  
 Beispiele:  
 0005 §218  
 0005 Abtreibung  
 0005 Gleichberechtigung
- Kategorie 0006 „Arbeit“  
 Beispiele:  
 0006 Arbeit  
 0006 Vollbeschäftigung  
 0006 Beschäftigung
- Kategorie 0007 „Recycling, Wiederverwertung“  
 Beispiele:  
 0007 Altstoffe  
 0007 Glaspfandsystem  
 0007 Müll
- Kategorie 0008 „Kriminalität“  
 Beispiele:  
 0008 Kriminalität  
 0008 Raub  
 0008 Rauschgift
- Kategorie 0009 „Kultur und Bildung“  
 Beispiele:  
 0009 Bildung  
 0009 Artistenschulen  
 0009 Kulturstätten
- Kategorie 0010 „Ruhe und Friede“  
 Beispiele:  
 0010- Ruhe (calm, quiet)  
 0010- friedlich

### 3.5 Reliabilität und Validität

Nach dem Aufbau des Diktionärs muß seine Reliabilität und Validität (Früh 1998), bezogen auf die Forschungsfrage, überprüft werden. Reliabilität bedeutet die Verlässlichkeit des Instruments, d.h. hier die Verlässlichkeit des Diktionärs: Wird eine Codierung mehrfach durchgeführt, muß sie immer zu gleichen Ergebnissen kommen. Diese Reliabilität ist bei der computerunterstützten Inhaltsanalyse immer gewährleistet, denn wenn man Codierungen mit dem Computer mehrfach für denselben Text mit demselben Diktionär durchführt, wird man immer wieder dieselben Ergebnisse/Codes bekommen.

Schwieriger ist die Sicherstellung der Validität. Validität bedeutet die Gültigkeit eines Instruments: Mißt ein Instrument tatsächlich das, was gemessen werden soll? Ist das theoretische Konstrukt angemessen erfaßt?



Zur Überprüfung der Validität und Verlässlichkeit des Diktions gibt es verschiedene Methoden:

- KWIC-Listen (TEXTPACK)  
Die Verwendung von KWIC-Listen wurde bereits bei der Erstellung des Diktions beschrieben. Mit Hilfe einer nach Kategorien sortierten KWIC-Liste für das gesamte Diktion läßt sich relativ schnell prüfen, ob die im Kontext gezeigten Wörtern „richtig“, d.h. ihrer Bedeutung entsprechend codiert wurden.
- Zurückschreiben der Codes in den Text (TEXTPACK)  
Ein weiteres Hilfsmittel der Überprüfung der Codierung bietet die Option, die Codierungen in den Text hinter (oder anstelle) des codierten Wort zu schreiben. Lesen des codierten Textes bietet eine einfache Kontrolle der Codierungen (siehe Tabelle 5). Beim Lesen des codierten Textes fällt auf, daß Familienbetreuung der Kategorie 4 („Soziale Beziehungen“) zugeordnet wird, richtig wäre aber die Kategorie 11 („Soziale Dienste“). Das Diktion muß entsprechend korrigiert werden.

In der Regel wird man die beiden oben beschriebenen Optionen auf eine handhabbare Auswahl der Texte (eine Textstichprobe) anwenden.

Tabelle 5: Zurückschreiben der Codierungen in den Text

```

***** T E X T P A C K      1 Aug 00   Routine -LIST-   *****
***** SENTENCE file: E:\TEMP\WAHL.SEN
08/02/

-ID1-  -ID2-  -ID3-  T e x t-----
000009          Sozialleistungen $$0001
000010          Für alle Arbeit $$0006 , soziale $$0001
                Dinge .
000014          Auf Familienbetreuung $$0004 .
000018          Kindergartensystem $$0011 , keine
                Arbeitslosigkeit $$0006 .
000020          Auf die Organisation für berufstätige
                $$0006 Mütter ( Kindertagesstätten $$0011
                ) ; geringe Arbeitslosenquote $$0006 ,
                geringe Kriminalität $$0008 .
000031          Soziale $$0001 Einrichtungen , Familien-
                $$0004 / Kinderbetreuung $$0011 ,
                Nachbarschaftshilfe $$0004
000038          Kinderhorte $$0011
000050          Auf die Disziplin , die Kinderbetreuung

```

- Liste der nicht-codierten Einheiten (TEXTPACK)  
Eine Liste aller nicht codierten Einheiten (in unserem Beispiel Antworten) hilft bei der Kontrolle der Vollständigkeit des Diktions. Eine Überprüfung muß zeigen, ob diese nicht-codierten Einheiten nicht codierbar sind oder ob im Diktion relevante Einträge fehlen. Im ersten Fall ist zu entscheiden, was mit den nicht-codierten Einheiten geschehen soll. Denkbar sind manuelles Nachcodieren oder Zuweisen eines fehlenden Wertes. Im zweiten Fall kann das Diktion entsprechend erweitert und die Codierung wiederholt werden.
- Liste der Einträge des Diktions, die zur Codierung führen (TEXTPACK)  
TEXTPACK bietet die Möglichkeit, eine Liste aller zur Codierung führenden Diktions-einträge zu erstellen. Die Liste kann für den gesamten Text oder pro Texteinheit erstellt

werden. Sie bietet die Möglichkeit zu überprüfen, ob die Einheiten richtig codiert werden, aber vor allem auch zu überprüfen, ob es Einträge im Diktionär gibt, die nie oder selten zur Codierung führen. Diese Einträge sind nicht genutzte Beispiele, die aus Gründen der besseren Übersicht und Handhabbarkeit aus dem Diktionär entfernt werden sollten.

- **Kategorienhäufigkeiten (TEXTPACK)**  
Am Ende des Codierprozesses wird eine Liste erstellt, die für jede Kategorie angibt, wieviele Textstellen ihr zugewiesen wurden. Hier kann man prüfen, ob die Kategorien angemessen definiert sind. Wurde z.B. eine Kategorie zu eng definiert (kaum Codierungen) oder eine andere zu breit (fast alle Textstellen gehören dazu)?
- **Manuelles Codieren einer Textstichprobe**  
Häufig wird auch der Weg gewählt, der bei der konventionellen Codierung üblich ist: das zweimalige Codieren einer Textstichprobe, d.h. ein Teil der Texte wird nicht nur automatisch, sondern auch von einem entsprechend geschulten Codierer codiert. Die Ergebnisse beider Codierungen werden verglichen und mit entsprechenden Reliabilitätskoeffizienten bewertet, z.B. mit einer unkorrigierter Verhältniszahl (Übereinstimmungen von Codierer 1 (= Computer) mit Codierer 2/Zahl der Codierungen), mit Kappa oder mit Scott's Pi (Merten 1995).

Nach diesen Prüfungen kann das Diktionär noch einmal verbessert werden, und danach kann die endgültige Codierung erfolgen.

## **4. Codierung und Analyse der Daten**

### **4.1 Codierung**

Im Gegensatz zur Inhaltsanalyse anderer Textarten, ist bei Antworten auf offene Fragen die Gliederung der Texte (die Textstruktur) in Text- und Codiereinheiten in der Regel bereits vorgegeben: Es empfiehlt sich meist, daß jede Antwort auf eine Frage eine eigene Texteinheit bildet. Codiereinheiten sind die einzelnen Antworten, d.h. während der Codierung wird jede Antwort einer oder mehreren Kategorien zugeordnet (bzw. jeder Antwort werden ein oder mehrere Codes zugewiesen).

Bei der Codierung wird jeder Einheit, sobald das entsprechende Wort im Diktionär identifiziert wird, der definierte Code zugewiesen. Dabei können für jeder Antwort mehrere Codes vergeben werden. Einer Antwort „Kindergartensysteme, keine Arbeitslosigkeit“ würden mit „0011“ und „0006“ codiert werden. TEXTPACK bietet verschiedene Möglichkeiten, diese Codes zu speichern: direkt im Text, d.h. der Code wird direkt in den Text geschrieben (siehe Validierung), oder als separate numerische Dateien für die weitere Analyse mit SPSS, SAS oder anderen Statistikprogrammen.

Die Codes können in der numerischen Datei als *Kategorienabfolge* oder als *Kategorienhäufigkeiten* gespeichert werden. *Kategorienabfolge* heißt, daß die Codes genau in der Reihenfolge gespeichert werden, in der sie im Text codiert werden. Im obigen Beispiel würden also in der Datei neben der Fragebogennummer die Codes „0011 0006“ gespeichert. Diese Kategorienreihenfolge ist z.B. bei einer Wortfeldanalyse von Bedeutung. Bei der Generierung von *Kategorienhäufigkeiten* wird pro Kategorie gespeichert, wie oft ein Code in einer Antwort codiert wurde. In unserem Beispiel würde die Ausgabe neben der Fragebogennummer für die Kategorien 1-5, 7-10 und 12 die Häufigkeit „0“, für die Kategorien 6 und 11 die Häufigkeit 1 enthalten.

In unserem Beispiel wurden durch das automatische Codieren 12 neue Variablen generiert. Jede Variable beschreibt eine unserer 12 Kategorien (z.B. die erste Variable definiert die Kategorie „Soziale Sicherheit allgemein“) und enthält zu jedem Befragten die Angabe, wie oft die entsprechende Kategorie von ihm genannt wurde.

## 4.2 Analyse der Daten

Hier soll nur ein kurzes Beispiel der weiteren Auswertung der Daten gegeben werden. Eine ausführlichere Analyse der Daten ist beschrieben in Mohler & Zuell (2001). In der Regel wird man die Ergebnisse der Codierungen der Antworten auf offene Fragen an die numerische Datei der Antworten auf die anderen Fragen des Fragebogens anfügen (als zusätzliche Variablen), um Auswertungen sowohl der geschlossenen wie auch der offenen Antworten zu machen.

TEXTPACK bietet optional die Erstellung einer SPSS-Datenbeschreibung als Syntax-File an, mit dem in SPSS ohne Aufwand eine Systemdatei aufgebaut werden kann. Diese Datei kann dann in SPSS mit Hilfe der Optionen „Data“, „Merge Files“ und „Add Variables“ an die Datei mit den Antworten auf die geschlossenen Fragen des Fragebogens angefügt werden.

Die nächste Entscheidung, die getroffen werden muß, ist die Behandlung der Kategorienhäufigkeiten. Die durch die computerunterstützte Inhaltsanalyse erstellten Variablen enthalten pro Befragten die Information, wie häufig diese Variable (=Kategorie) genannt wurde. Hier ist nun die Frage zu beantworten, ob Häufigkeiten eine Bedeutung haben (was z.B. bei der Analyse von Rhetorik von Wichtigkeit ist) oder ob nur die Information Kategorie „trifft zu“/„trifft nicht zu“ von Interesse ist. Da für unsere Fragestellung die Häufigkeit der Nennung einer Kategorie bei einem Befragten keine Rolle spielt, haben wir uns entschieden, die Variablen zu recodieren, so daß wir in der weiteren Analyse nur die Information „trifft zu“/„trifft nicht zu“ verwenden (0/1).

Tabelle 6: Stolz auf "Arbeit" (Arbeitsplatzsicherheit) in Ost- und Westdeutschland abhängig von der gewählten Partei

Split Ost/West			Arbeit		Gesamt
			Kategorie nicht genannt	darauf kann man Stolz sein	
<b>West</b>	Zweitstimme	CDU/CSU	56	13	69
		SPD	68	25	93
		FDP	4	3	7
		Grüne	40	9	49
		PDS	13		13
		Republikaner	3		3
		andere	3		3
	Gesamt		187	50	237
<b>Ost</b>	Zweitstimme	CDU/CSU	89	83	172
		SPD	83	108	191
		FDP	10	10	20
		Grüne	14	12	26
		PDS	58	74	132
		Republikaner		1	1
		andere	4	2	6
	Gesamt		258	290	548

In Tabelle 6 ist als Beispiel eine Auswertung<sup>2</sup> gezeigt. Alle Befragte, die angaben, daß es nichts gibt, auf das man stolz sein konnte, wurden in der Tabelle ebenso ausgeschlossen wie die 281 Befragten, die die Frage, welcher Partei sie bei der letzten Wahl ihre Zweitstimme gegeben haben, nicht beantwortet haben. Überprüft werden soll, ob es zwischen Befragten in Ost- und Westdeutschland in Abhängigkeit vom Wahlverhalten (Vergabe der Zweitstimme) Unterschiede bei der Nennung der Kategorie „Arbeit“ (Arbeitsplatzsicherheit) als Objekt des Stolzes gibt. Die Tabelle zeigt, daß nur sehr wenige westdeutsche Befragte, die Arbeitsplatzsicherheit in der ehemaligen DDR als etwas betrachten, auf das man stolz sein kann. Dagegen scheint das Thema für ostdeutsche Befragte ein wichtiger Aspekt zu sein. Das unterschiedliche Wahlverhalten der Befragten hat wenig Bedeutung für die Relevanz des Themas „Arbeit“, obwohl es für ostdeutsche SPD- und PDS-Wähler ein wichtigerer Aspekt ist als für Wähler anderer Parteien.

Mit den vorliegenden Daten kann z.B. eine multidimensionale Analyse durchgeführt werden, um Informationen über den Zusammenhang der Kategorien zu erhalten. Diese Analyse wird für Ost- und Westdeutschland getrennt durchgeführt (Tabelle 7 und 8). Bei Antworten auf offene Fragen bietet sich die MDS an, weil sie ein dem Datenniveau und der Textbasis angemessenes, strukturabbildendes Verfahren ist. Andere Verfahren wären etwa Clusteranalyse oder Korrespondenzanalyse. Häufigkeitsauszählungen und Kreuztabellen können nur ein Hilfsmittel auf dem Weg zu dieser Analyse sein.

In beiden Fällen zeigt sich eine zentrale Clusterung mit einigen peripheren Kategorien. Bei den westdeutschen Befragten zeigt sich eine Gruppe von Kategorien/Themen im Mittelpunkt umgeben von einigen Kategorien als Ring. Im Zentrum liegen die Kategorien „Kultur & Bildung“, „Kosten & Preise“, „Kriminalität“, „Ruhe & Friede“, „Recycling“, „Medizinische Versorgung“ und „Frauenrechte & Abtreibung“. Wenn westdeutsche Befragte mehr als einen Bereich nennen, der zum Stolz Anlaß gibt, dann sind dies Kategorien im Zentrum. Die Kategorien „Arbeit“, „Sport“, „Soziale Dienste“, „Soziale Sicherheit allgemein“ und „Soziale Beziehungen“ werden dagegen in der Regel alleine, d.h. ohne eine zweite Kategorie benannt.

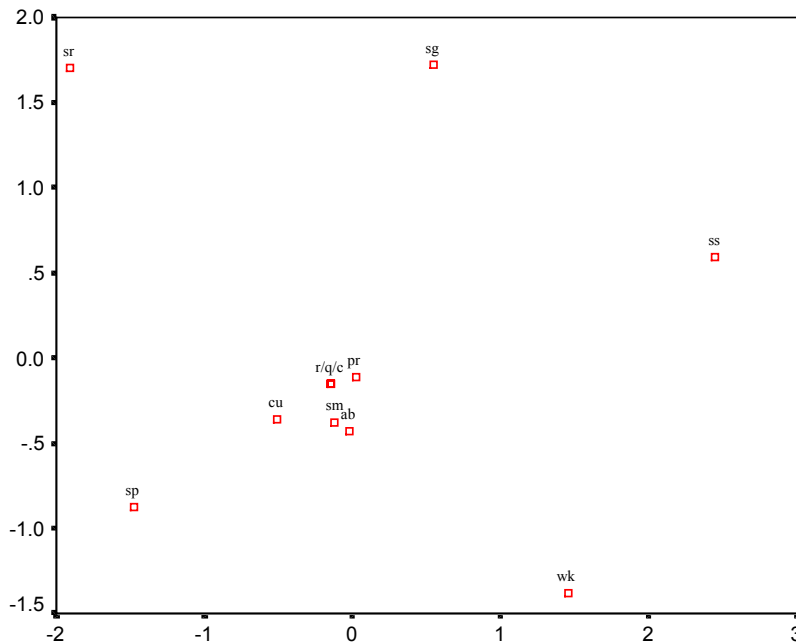
In Ostdeutschland dagegen sind 8 der 12 Kategorien eng miteinander verbunden, d.h. sie werden jeweils gemeinsam genannt. „Soziale Dienste“ und „Kosten & Preise“ bilden dabei eine eigene kleine Untergruppe. Dies bedeutet, daß für ostdeutsche Befragte diese Kategorien eng in Beziehung zueinander stehen. Die Kategorien „Arbeit“, „Soziale Sicherheit allgemein“, „Kultur & Bildung“ und „Soziale Beziehungen“ werden dagegen klar getrennt. Auffällig ist auch, daß die Kategorie „Kultur & Bildung“ für westdeutsche Befragte ein wichtiger Aspekt zu sein scheint, für ostdeutsche dagegen nicht. Dies ist ein Hinweis auf die unterschiedlichen Perspektiven direkt Betroffener und Außenstehender.

Auf eine ausführliche Interpretation der Ergebnisse und zusätzliche Analysen soll an dieser Stelle nicht weiter eingegangen werden. Das hier gegebene Beispiel soll nur zeigen, wie die Daten einer Inhaltsanalyse mit komplexen statistischen Verfahren weiteranalysiert werden können.

---

<sup>2</sup> Alle statistischen Analysen wurden mit SPSS 9 durchgeführt.

Tabelle 7: Ergebnisse der Multidimensionale Skalierung (euklidische Distanz) der westdeutschen Befragten

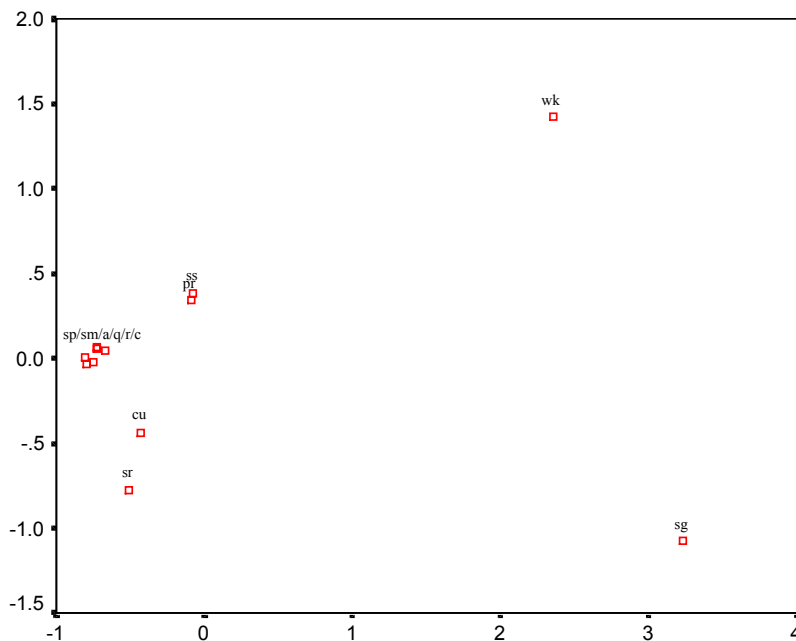


Legende:

pr = Kosten & Preise	sp = Sport	sr = Soziale Beziehungen	ab = Frauenrechte & Abtreibung
wk = Arbeit	r = Recycling	c = Kriminalität	cu = Kultur & Bildung
q = Ruhe & Friede	ss = Soziale Dienste	sm = Medizinische Versorgung	sg = Soziale Sicherheit allgemein

Achtung: Die Kategorien Kosten & Preise, Ruhe & Friede und Recycling überlappen im Diagramm.

Tabelle 8: Ergebnisse der Multidimensionale Skalierung (euklidische Distanz) der ostdeutschen Befragten



Legende:

pr = Kosten & Preise	sp = Sport	sr = Soziale Beziehungen	ab = Frauenrechte & Abtreibung
wk = Arbeit	r = Recycling	c = Kriminalität	cu = Kultur & Bildung
q = Ruhe und Friede	ss = Soziale Dienste	sm = Medizinische Versorgung	sg = Soziale Sicherheit allgemein

## **5. Weiterführende Informationen zur Inhaltsanalyse**

Die auf Diktionären basierende computerunterstützte Inhaltsanalyse wird bei Stone (1966) beschrieben. Neben dem hier erläuterten Ansatz der quantitativen computerunterstützten Inhaltsanalyse, der auf einem vom Anwender entwickelten Diktionär basiert, gibt es verschiedene andere Ansätze, auf die in diesem Papier nicht eingegangen werden kann. Eine Darstellung der verschiedenen Ansätze der computerunterstützten Inhaltsanalyse finden sich u.a. in Alexa (1997), Popping (2000) und in Roberts (1997). Ansätze der qualitativen Inhaltsanalyse sind z.B. in Flick et al. (2000) zusammengefaßt.

Ausführliche und sehr praxisbezogene Beschreibungen zur Methode der (konventionellen) Inhaltsanalyse bieten Früh (1998) und Merten (1995) an. Beide gehen allerdings nur sehr kurz auf die computerunterstützte Inhaltsanalyse ein.

## **Literatur**

Alexa, M. (1997): Computer-assisted Text Analysis Methodology in the Social Sciences. Arbeitsbericht 97/07. Mannheim: ZUMA.

Flick, U., von Kardorff, E., Steinke, I. (Hrgs.) (2000): Qualitative Sozialforschung. Ein Handbuch. Rowohlt.

Früh, W., (1998): Inhaltsanalyse - Theorie und Praxis. Konstanz: UVK Medien.

Merten, K. (1995): Inhaltsanalyse: Einführung in Theorie, Methode und Praxis. Opladen: Westdeutscher Verlag.

Mohler, P.Ph., Zuell, C. (1995): TEXTPACK User's Guide. Mannheim: ZUMA.

Mohler, P.Ph., Zuell, C. (2001): Applied Text Theory: Quantitative Analysis of Answers to Open-Ended Question. In West, M. D. (ed): Applications of Computer Content Analysis. Westport: Aplex Publishing, S. 1-16.

Popping, R. (2000): Computer-assisted Text Analysis. London: Sage.

Roberts, C.W. (Hrg.) (1997): Text Analysis for the Social Sciences: Methods for Drawing Statistical Inferences from Texts and Transcripts. Mahwah, N.J.: Lawrence Erlbaum Assoc.

Stone, P.J., Dunphy, D. C., Smith, M.S., Ogilvie, D. M. (1966). The General Inquirer: A Computer Approach to Content Analysis. Cambridge: The M. I. T. Press.