

## CHAPTER 14

# A Practical Guide

*This final chapter discusses three starting points for content analyses. For each, it recommends procedural steps, raises issues that might come up during the research, notes the junctures at which content analysts need to make decisions, and suggests what they need to take into consideration in making those decisions. It gives an overview of the entire content analysis process, from conceptualizing the research questions that content analysts are called on to answer to reporting their results, providing ample references to related material in the foregoing chapters.*

**I**n the preceding chapters, I have introduced the concepts involved in content analysis one by one and have suggested solutions to conceptual and methodological problems that content analysts need to address. In this chapter, I rearticulate these concepts with practice in mind, so that readers who have particular research problems that content analysis might solve will get an idea of what they can or need to do and why.

Like most social research, content analysis involves four kinds of activities:

- Designing an analysis (section 14.1)
- Writing a research proposal (section 14.2)
- Applying the research design (section 14.3)
- Narrating the results (section 14.4)

These activities are neither mutually exclusive nor entirely sequential. Surely, all inquiries start with some conceptualization of the research process, but as content analysis requires a great deal of preparatory effort, the design phase may well become part of a research proposal. Researchers who have their own resources, including students working on their theses, may not need to write formal research proposals, but they are nevertheless well-advised to have such proposals in mind. Often, the relationship between the design of a content analysis and the application of that design is circular. A seemingly perfect design may reveal flaws in its

application that bring the researcher back to the drawing board to sample more or different kinds of data, reconceptualize the data language, improve the coding instructions, and even radically change the approach taken initially. This is the hermeneutic circle in scientific inquiry.

## Designing an Analysis 14.1

Any design proposes something that would not come into being without guided efforts—here, a procedure for moving from potentially available observations, texts, sounds, and images to the narrated answers to a research question (see Figure 4.2). A research design consists of the detailed specifications that guide the handling of data and make the research reproducible and critically examinable at a later point in time. Aside from getting involved in available text, the development of a research design is the most intellectually challenging part of a content analysis. In the course of designing a study, analysts clarify their research interests, learn to respect their own readings as different from those of others, explore the questions they would like to see answered, delve into available literature for insights about the contexts of their analysis, and play with analytical possibilities—until step-by-step specifications emerge that promise to bring the analysis to a worthwhile conclusion (see Chapter 4).

Researchers may enter content analysis from different starting points. The possible starting points may not be equally desirable, but circumstances for research rarely are. Below, I discuss content analyses in terms of three points of entry:

- *Text-driven* content analyses (14.1.1) are motivated by the availability of texts rich enough to stimulate the analysts' interests in them. As the research questions emerge, as analysts are becoming involved with such texts, text-driven analyses are also called "fishing expeditions."
- *Problem-driven* content analyses (14.1.2) are motivated by epistemic questions about currently inaccessible phenomena, events, or processes that the analysts believe texts are able to answer. Analysts start from research questions and proceed to find analytical paths from the choice of suitable texts to their answers.
- *Method-driven* content analyses (14.1.3) are motivated by the analysts' desire to apply known analytical procedures to areas previously explored by other means.

### 14.1.1 Text-Driven Analyses

A text-driven analysis starts with a body of text, as noted above: an interesting set of personal letters, a collection of taped interviews, the diary of a famous person, a compilation of comic books, transcripts of naturally occurring conversations, conference proceedings (a collection of presented papers), a significant period of publications (newspapers, professional journals, movies), election campaign speeches, news accounts of a particular crime, notes made during anthropological

fieldwork, a collection of family photographs, advertisements in magazines published in different countries, reports to the shareholders of a corporation, library listings of a particular institution, articles mentioning a particular drug, telephone books, and so on.

Without an explicit research question in mind, researchers typically start by familiarizing themselves with the chosen body of texts. They may begin with “housekeeping” chores—cataloging the texts, unitizing the body of text into “packages” that can be handled more or less independent of each other. When the researchers have a sense of how many texts there are, they may want to explore apparent intertextualities: quotes, references, overlaps, rearticulations, elaborations, and sequential orderings. They may also note how the texts reproduce, respond to, or elaborate on each other. Intertextualities form dependency networks among texts, which in turn suggest possible paths for reading the given body of texts—for example, along lines of narrative coherences, genres, or historical sequences, or according to how various sources respond to each other. Constructions of intertextualities amount to ways of reading, but—to be clear—these always are an analyst’s readings.

Next comes a careful reading of the texts for the purpose of summarizing what they collectively mean to the analyst, what they denote, connote, or suggest, or how they are or could be used as a whole. If the body of texts is so large that it becomes difficult for a single individual to keep track of all the relevant details while reading, content analysts may take advantage of a variety of computer aids, including qualitative data analysis (QDA) software such as QSR’s NVivo, R’s RQDA, or word frequencies and KWIC lists with software such as Concordance, VBPro, WordStat, or TextQuest. Even the search functions of ordinary word processing programs, such as Microsoft Word, may be used for simple content analyses. ATLAS.ti is especially well suited to the creation of networks of texts. The advantage of such analytical aids lies in the assurances they can provide that text explorations are systematic, effectively countering the natural tendency of humans to read and recall selectively. Text-driven analyses are also called *interpretive* or *qualitative*, as opposed to *quantitative* (see Chapter 4, section 4.1.2), discussed as QDA in Chapter 11 with reference to computer aids (see section 11.5); such analyses are indebted to the tradition of grounded theory beginning with Glaser and Strauss (1967).

Even when aided by QDA software, it is important to note, such text explorations are essentially limited to a single analyst’s conceptions and ability to read. Such software is convenient, but it does not assure the replicability of the process.

When the volume of texts exceeds individual abilities as well, when teamwork becomes essential, the problems of coordinating the work of several analysts begin to dominate the research effort. It then becomes important that the analysts working on a project agree on the terms of the analysis; in addition, as differences in their individual readings will inevitably surface, the analysts need to render these readings in comparable categories. This will bring the analysis in line with the procedures described in Chapter 4—whether the analysis is called qualitative or quantitative. I continue this thread below, in section 14.1.2.

Problems of coordination among content analysts on a large project also challenge the analysts to overcome the temptation to take their own readings as the only ones that count, which often coincides with the belief that “content” is contained in text, an inherent quality of text, and that everyone ought to know what *it* is (see Chapter 2). Obviously, texts can be read variously and mean different things to different people. Recognizing this leads analysts to listen, while reading texts, to the voices of other readers—the texts’ writers, audiences, and users—to understand what the texts mean to them, how language is implicated in their interpretation, the roles the texts do or could play in the lives of their users, and whose voice is represented, who listens, who responds, and who is silenced in these texts. Unless they consider alternative readings, content analysts are limited to analyzing their own understanding and nothing outside it. By reading with alternative voices in mind (see Krippendorff, 2006), analysts begin to expand their horizons, to get a feeling for the context they need to construct to make room for diverse people associated with the texts and for the institutions that govern the texts’ antecedent conditions, concurrent interpretations, and practical consequences. As this context becomes conceptually clear, so do the questions that analysts feel they could answer. Figure 4.3, which depicts this text-driven entry to content analysis, suggests that such questions arise not only from the virtual voices of other readers or users but also from what the authors of pertinent literature on the context say.

Qualitative content analysts typically stop at their own interpretations of texts, however, and the recommendation I have just made is intended to undermine the very notion that texts can drive an analysis. Contrary to what lawyers are fond of saying, documents never speak for themselves—interpretations are always made by intelligent readers. And texts inevitably have multiple meanings. Although the seeming objectivity of computer searches tends to hide this fact, the queries that inform such searches are always formulated by the very analysts who also interpret the search results. When analysts acknowledge their own conceptual contributions and, by implication, the possibility of diverse readings, especially when they have to train coders to comply with written recording/coding instructions, their analyses can no longer be conceived as text driven. This brings us to the next, perhaps more preferable, way of entering a content analysis.

### 14.1.2 Problem-Driven Analyses

A problem-driven analysis derives from epistemic questions, from a desire to know something currently inaccessible and the belief that a systematic reading of potentially available texts and other data could provide answers. Content analysts who start at this point tend to be involved in real-world problems: psychoanalysts seeking to diagnose the pathology of a patient; lawyers trying to generate evidence that will support or refute an accusation of plagiarism; historians hoping to clarify how a historical event unfolded; literary or forensic scholars aspiring to know who wrote an unsigned document; educators attempting to predict the readability of textbooks; mass-media researchers intending to substantiate claims of undesirable effects on civic society of TV election coverage; propaganda analysts looking for

military intelligence hidden in enemy domestic broadcasts; rhetoricians desiring to measure how much civic knowledge candidates for political office bring to a debate. All of these are epistemic problems—that is, they are problems of not knowing something deemed significant. Content analysts must convert such problems into research questions, which they then attempt to answer through a purposive examination of texts. In Chapter 3, I discuss the different kinds of inferences that analysts may draw from texts; here, the focus is on the steps that analysts may want to take to get to those inferences:

- Formulating research questions
- Ascertaining stable correlations
- Locating relevant texts
- Defining and identifying relevant units in texts
- Sampling these units of texts
- Developing coding categories and recording instructions
- Selecting an analytical procedure
- Adopting standards
- Allocating resources

#### *14.1.2.1 Formulating Research Questions*

Formulating research questions is by far the most important conceptual task that analysts face, for it is the key to a successful research design. As noted in Chapter 2, the research questions posed in content analysis have the following characteristics:

- They concern currently unobserved phenomena in the problematized context of available texts.
- They entail several possible answers.
- They provide for at least two ways of selecting from among these answers—if not in practice, then at least in principle.

The first of these defining features rearticulates the aim of content analysis: to make abductive inferences from texts to phenomena outside those texts. In Chapter 2 (section 2.4), I discuss how abduction is distinguished from induction (moving from particulars to generalizations) and deduction (moving from generalizations to particulars) as a form of reasoning that moves from particular texts, through context-sensitive explanations of these texts, to particular answers to research questions (i.e., from particulars to particulars). Such answers need to solve a problem in the widest possible sense, which may range from solving an important mystery to informing a decision to act.

The second characteristic above demands of a research question that it entail several conceivable answers—neither an open field in which anything can happen nor a single answer that the analyst intends to prove.

The third feature suggests that it would not be sufficient simply to answer a research question, even for good reasons. The researcher should think of at least

one other way to answer that question, a way not involving content analysis that could validate the answer that the content analysis suggests, at least in principle: additional observations, correlations with measures already known to be valid, even observations of improved success in acting on the information the content analysis provided (see Chapter 13). A content analysis should be validatable in principle, and this entails an alternative method or independently available evidence.

Novices in content analysis need to understand that not all questions qualify as research questions. For example, a question concerning whether a computer program is suitable for analyzing a body of text is a question about the properties of the software relative to available texts—it does not address anything in the context of the analyzed texts. The question of whether one can measure a desired quality of texts does not qualify either, because it concerns an analyst's ability and is answerable by that analyst's doing it. The latter has nothing to do with the chosen context of analysis. Nor are questions concerning how often an author uses a certain expression appropriate research questions. Counting is an operation performed on a body of text. Its result says nothing other than that someone has counted something. This argument applies to all computational abstractions: type/token ratios, vocabulary sizes, and various correlations within a text. Unfortunately, the history of content analysis is full of examples of researchers who have merely declared counts to be indices of phenomena, usually of social or political significance, without spelling out how their claims could be validated. For example, a count of the number of violent acts on television means nothing unless one has at least a hunch that high numbers of such acts are bad for something one wants to preserve. More pernicious are pseudoquestions that exclude alternative answers. In the 1930s, many journalists quantified categories of content as a way of "objectifying" preconceived public concerns. Content analysts must not confuse proving one's point in a public debate with pursuing a research question that has several possible answers.

In effect, the third definitional feature of research questions calls for their answers to be validatable in principle (see Chapters 2 and 13)—in principle because many content analytic situations preclude validation in practice. For example, history always happened in the past. It no longer is observable, although it arguably was at one time. The traces of historical events that do survive are nothing more than indirect indicators of the events themselves. Texts are no exception. Inferences about historical events are best correlated with other traces left behind. To find relevant traces, assess their trustworthiness, and relate them to appropriate dimensions of a content analysis, analysts require considerable insight into the context of the analysis. Psychotherapeutic diagnosis provides a different example. No therapist can enter a patient's mind, so the first definitional feature of research questions is satisfied. The American Psychiatric Association's *Diagnostic and Statistical Manual of Mental Disorders* (2000) lists all legitimate mental disorders for therapists to choose from. One may disagree with this list, but it does satisfy the second feature. The diagnosis, an informed selection among known mental disorders, might be confirmed by other therapists, validated by independent tests, or vindicated by a successful treatment. Similarly, the answer to the question of whether an author committed plagiarism may not be known for sure without

witnesses or admission by the accused. The latter could validate or invalidate the conclusion derived from a systematic comparison of two texts. Some of the more traditional questions that content analysts have answered concerned what different audiences could learn from exposure to particular mass-media messages and the political leanings of newspaper editors. To validate inferences from these kinds of media messages, content analysts have compared their results with survey results, with information from interviews with experts, and with focus group data—all sources of validation on which social scientists rely heavily. Content analysts should be specific about what could validate or invalidate their results, even if no (in)validation effort is ever undertaken, because in doing so they tie content analysis to a multiply confirmable reality.

Well-formulated research questions not only guide a research design, they also constitute one-third of a content analyst's world construction, the framework within which the analysis is to be undertaken. Research questions concern the uncertain or variable part of that world (see Figure 2.1). The analysts' next two steps involve the stable parts of their world.

#### 14.1.2.2 *Ascertaining Stable Correlations (With the Research Questions)*

Content analysis answers research questions by analyzing texts, which are understood quite generally to include images, sound, websites, symbolic events, even numerical data, provided they mean something in the chosen context. The abductive inferences that content analysis entails presuppose some knowledge on the analysts' part of how the research questions relate to available texts. This knowledge need not be, and typically is not, exclusively linguistic or semantic. Content analysts are rarely interested in what is said literally, by dictionary definition or according to a standard reader, if such a person exists. Content analysts are as interested in what is not said as they are in what is said—that is, they are interested in what texts reveal about phenomena not spoken of, such as ideological commitments or ethnic prejudices that are manifest in influences, consequences, and uses that may well go unrecognized by individual readers. Inferences of the latter kinds tend to rely on statistical knowledge and call for the use of complex analytical instruments. In Chapter 2, I describe these connections as stable correlations, stable or enduring because only if they can be assumed to remain invariant, at least up to the conclusion of an analysis, can they justify the inferences that content analysts are asked to make. Questions whose answers are not correlated with anything observable, readable, or accomplishable with texts cannot be answered.

Traditionally, content analysts have focused on *linguistic references*, *expressions* of attitudes, and *evaluations*. These assume a one-to-one correlation between textual units and the phenomena referred to, expressed, articulated, or, in a naive sense, “contained” in them. More recently, content analysts have relied on texts as *statistical correlates* of the phenomena of interest, using the wear and tear shown by library books to assess the popularity of certain topics, for example, or relying on the correlation between expressed public concerns and voting. Typically, such

correlates stem from other methods of inquiry—such as public opinion research, media effects studies, perception experiments, and theories of cultural cognition—that often are concerned with frequencies, contingencies, and variances. Texts may be seen also as *by-products* of the phenomena of interest, such as when researchers use mass-media coverage to infer how the mass-media industry is organized; as *causes*, such as when researchers attempt to infer audience perceptions or media-induced anxieties; as *consequences*, such as when researchers analyze medical records to determine the population characteristics of patients; or as *instrumental*, such as when researchers take texts as evidence of manipulation efforts by their authors, as in political or public health campaigns. Webb, Campbell, Schwartz, and Sechrest (1966) add to this list of possible connections *physical traces* that the phenomena of interest leave behind and *actuarial records* that institutions maintain for reasons other than analysts' interest in them. Dibble (1963) recognized the latter as well, and George (1959a) has reported the use of complex sociopolitical networks of stable correlations to answer research questions (see Chapter 9 and Figure 9.2). Artificial intelligence models account for still other stabilities in the world of content analysts, enabling them to obtain answers from sparse textual evidence.

The analysts' task at this step of a research project is to ascertain a reliable network of these correlations, correlations that researchers can rely on to be stable and general (invariant over time and in various situations), certain (able to determine or be determined), and selective (able to narrow the set of possible answers to a research question). Content analysts may utilize all kinds of sources for this knowledge. Past empirical research about text-context relationships is one such source. In Chapter 4 (section 4.2), I suggest several designs for the generation of empirical knowledge in preparation for a content analysis. Available theories and models of how the phenomena of interest are communicated within a social system can also serve as sources of needed evidence. Chapter 9 gives examples of analytical constructs that are as simple as the aforementioned one-to-one correlations and as complex as a model of how a government controls public opinion through its publications, such as in preparatory propaganda efforts. Causality is not the only path, as suggested by the above list of correlations, nor do correlations need to be direct. Content analysts often lament the absence of general theories and the naive simplicity of theories that claim universality. Indeed, literature on the contexts of particular texts, another important source of knowledge of prevailing stabilities, often includes situation-specific descriptions, temporally limited “mini-theories,” and highly qualified propositions, all of which could well support a content analyst's efforts. These if-then propositions may be derived from ordinary readings, known stereotypical reactions, widely used folk sayings, metaphors, colloquialisms, and so on. The hope is that such a network of propositions contains a path that connects available texts to the needed answers. A well-substantiated path is all that content analysts need to create.

In assuming such a stable network of correlations, analysts may need to keep track of the conditions under which these stabilities are warranted and when they become unreliable. For example, correlations that have been found to hold for



undergraduate subjects may not be generalizable to other populations, at least not without qualifications. In crisis situations, organizational rules may break down or be replaced by others. The meanings of verbal expressions may change over time and/or become variable from one social situation to another, or from one culture to another. Some content analysts are tempted to assume linguistic universality or to assume that correlations once found do not change under conditions other than those studied (see Figure 2.1), but assuming such correlations to be stable when they are not can seriously mislead content analysts.

Obviously, knowledge of the needed correlations is informed by how a context is defined. Political scientists look for correlations that differ from those that psychiatrists are able to consider. Sociologists and communication researchers approach content analyses with different constraints in mind, and the worlds that they respectively construct for given texts may well be incommensurate with one another.

A network of stable correlations constitutes the second third of the content analyst's world construction (see Figure 4.2). Its purpose is to channel, almost in an information theoretical sense, the diversity encountered in texts to the possible answers to a research question. The analysts' next step is one that is often mentioned as the starting point for content analysis: locating relevant texts.

#### 14.1.2.3 Locating Relevant Texts

In content analysis, texts inform analysts' questions and so must be sampled from populations of texts that can be informative in this sense. A text is relevant if there is evidence for or an assumption of stable correlations between that text and answers to the research question. By backtracking along the path of the intended inferences, moving from the phenomena of interest along the stable correlations to potentially available texts, content analysts can justify the relevance of a population of texts to given research questions.

As noted above, traditionally content analysts have made their analytical efforts easy by assuming one-to-one relationships between textual units and the phenomena of interest, what they are assumed to refer to, express, or "contain." With this assumption, selecting (including reading and counting) texts virtually substitutes for selecting (including observing and enumerating) the phenomena addressed by the research questions. This hides the usually complex roles that texts play in social situations. For example, to infer variations in achievement motives over various periods of a culture, McClelland (1958) searched for messages in which achievements were *created*, *negotiated*, or *celebrated*. This led him not only to popular literature, biographies, and expressions in art but also to images on Greek vases and postage stamps. Content analysts have a long history of analyzing the newspapers read by political elites in various countries to infer the politics of the citizens of those countries. This choice is grounded in the assumption that political agendas are set and public debates are spearheaded by certain leading newspapers, so-called prestige papers, rather than by local ones, which are more likely to reproduce what the prestige papers print and are, hence, less informative. To reveal an author's

suspected racial or ethnic prejudices, content analysts may have to sample from the writings of that author not intended for publication, such as personal diaries, letters to close acquaintances, or texts written specifically for the author's ethnic in-group. When a source of texts has a stake in the outcome of the analysis, the content analysts need to consider what that source knows about how it might be analyzed and focus instead on textual characteristics that the source is not aware of or cannot control. This rules out instrumental aspects of communication.

Although the logic of such choices is pretty clear, it is not easy to be more specific about how content analysts go about deciding on the informativeness of the texts they propose to analyze. Reading a small sample is a good start. Examining headlines or abstracts for clues to the relevance of texts is a common practice. Pursuing citation networks to the key publications is a strategy familiar as snowball sampling (see Chapter 6, section 6.2.6). Alleged expertise could also lead analysts to suitable populations, provided such attributions can be trusted. The reputation of a publication is another criterion for locating relevant texts. Often, however, content analysts have to be satisfied with what is made available to them. Propaganda analysts during wartime and analysts working for institutes that monitor international agreements have to start with what they can intercept. Scholarship on literary figures is limited to what is written by and about them and their times. Conversation analysts can record only with permission, which excludes many privileged conversations (and can introduce unwanted biases).

The internet as well as large, full-text electronic databases and digital libraries have vastly expanded the availability of content-analyzable texts. Browsers, text search engines, and computer-aided text analysis tools (see Chapter 11, section 11.3) can locate relevant texts in stages. Starting perhaps with vague hunches about what is relevant, text mining may begin a search with queries that cast a deliberately wide net over all conceivable texts just to get a sense of how much is there. Making good use of abduction, analysts typically develop increasingly detailed explanations of the body of texts they have so far scanned, become clearer about how available text may be correlated with the research question, and iteratively narrow the search. During such explorations, content analysts develop conceptions of the available texts and how to analyze them while simultaneously reducing the sample of texts to a manageable size. Search engines typically are severely limited in what they can identify (see Chapter 11, section 11.3). There almost always remains a considerable gap between what they retrieve and what is relevant (see the discussion of semantic validity in Chapter 13, section 13.2.2). Nevertheless, techniques for sampling from electronic databases, Web pages, and on-line exchanges are improving, promising content analysts an increasingly rich source of textual data.

Increasingly, content analysis is applied to texts obtained from internet discussion groups, e-mail exchanges, websites, blogs, Facebook, and Twitter. This can pose ethical issues. True, most content analyses rely on public documents, newspapers, books, radio broadcasts, corporate accounts to their shareholders, or historical documents in the public domain, many of which have been generated by institutions or authors long deceased. In content analysis, ethical issues have been largely delegated to those generating textual matter. Therapeutic interviews, for example,

typically are recorded by therapists, whose professional code of conduct prevents them from revealing the identities of their clients when publishing their findings. Although the distinction between public and private spheres is not always clear, I distinguish two kinds of internet data of interest to content analysts. The analysis of websites, blogs, court records, and documents that are meant to be public is ethically neutral. However, content analysts can easily and unintentionally violate the privacy of individuals whose writing is not intended to be public, are unaware that their texts are being analyzed for unintended purposes, and would take offense at such a violation of their privacy. On the internet, even assigning pseudonyms rather than real names to individual voices may no longer provide sufficient protection of the identities of writers or speakers, as powerful search engines may be able to find any given quote's source. Eysenbach and Till (2001) have distinguished three methods of gathering internet data from discussion groups, which may be described briefly as unobtrusive, with consent, and through explicit participation in deliberations; they discuss the methodological pitfalls of all three. Sixsmith and Murray (2001) go further in discussing the ethical issues associated with analyzing internet posts and archives in terms of consent, privacy, anonymity, and the issue of interpreting the meanings of others' voices.

Despite these constraints, with well-formulated research questions, a good sense of the network of correlations operating in the chosen context, and respect for the ethical issues of using the voices of others in ways they did not intend, content analysts can collect a wide range of data to explore. The analysts' decision on the population of texts to be considered completes the construction of the world in which the content analysis can proceed (see Chapter 2). We now turn to the components of content analysis (see Figure 4.2).

#### 14.1.2.4 Defining and Identifying Relevant Units in Texts

One way to make a content analysis of a large volume of text manageable is to break it into smaller units and deal with each separately. In Chapter 5, three kinds of units are distinguished according to the functions they serve within the content analytical process: *sampling units* are mutually exclusive units of text that are selectively included in an analysis (see Chapter 6); *recording units* are also mutually exclusive, either equal to or contained in the sampling units, but separately described, coded, or recorded in the terms of a data language; and *context units* set limits on the amount of text to be consulted in determining what a recording unit means (see Chapter 7). Chapter 5 also mentions *units of enumeration*, which usually coincide with recording units, sometimes in the form of numerical measurements: column inches, type sizes, picture sizes, ratios of different kinds of words, and other numerical scales.

In text mining, units can be defined by proximity operators (see Chapter 11, section 11.3.2): a delineated stretch of text, a document, an article, or a paragraph that contains a match with the query. Search results may serve as sampling units or as recording units. It is conceivable that text searches provide answers to analysts' research questions directly, but this is rare.

The definitions of units of analysis have important implications. When the units are mutually exclusive, counting them leads to comparable frequencies; when they overlap, it does not. Separating units also severs all relations among them, omitting information that resides between neighboring words or phrases. For example, taking single words as units disregards their roles in sentences, so that their syntactical meanings are lost; units in the form of sentences omit the roles that sentences play in paragraphs, thought sequences, the points made in longer arguments, and so on. Thus *unitizing a text is justifiable only if the relationships between units do not inform the research question*. The above-mentioned context units are intended to preserve at least some of the information that surrounds the recording units. Generally, if units are too small (such as words or short expressions), semantic validity suffers and a content analysis tends to become shallow. If units are too large (e.g., whole documents, Web pages, books, TV shows), the content analysis tends to become unreliable (see Chapter 12, section 12.1).

#### 14.1.2.5 Sampling the Texts

If the population of relevant texts is too large, content analysts may select representative samples of these texts. A heuristic approach to sampling is to start with any arguably unbiased sample of text, analyze it for how well it answers the research questions, and, if it fails to meet acceptable standards, either continue to sample until the questions are answered with reasonable certainty or declare the process hopeless. The latter may signal the need for a redesign of the content analysis.

Chapter 6 suggests suitable sampling strategies. However, because texts as well as the targets of content analyses are about phenomena outside or surrounding the texts, sampling in content analysis differs from sampling in other research techniques—for example, sampling of individuals for public opinion surveys. In content analysis, researchers need to sample texts with two populations in mind: the “population” phenomena that correlate with and hence lead to answers to a research question and the population of texts that represents these phenomena. In sampling the texts for a content analysis, researchers must give both phenomena a fair chance of contributing to the answer to the research question.

Chapter 6 also addresses the problem of texts that come into analysts’ hands for reasons unrelated to the research questions (see also the discussion of text-driven analyses in section 14.1.1). If texts are made available rather than purposefully sampled, their representativeness cannot be assured. They may be biased on account of their sources’ selectivity. For example, historical documents survive for a variety of reasons that are usually unrelated to why they are of interest. Politicians have good reasons to hide embarrassing information from public view, and television news is not about what happens in the world but rather represents what the mass-media institutions deem newsworthy and what they can fit into available programming space. When the population of available texts is small, content analysts may not have the luxury of sampling and so face the problem of rectifying the sampling biases inherent in these texts; see sampling validity (b) in Chapter 13, section 13.2.1.

#### 14.1.2.6 Developing Coding Categories and Recording Instructions

As I have noted above, when the volume of text exceeds a single researcher's analytical capabilities and analysts must therefore work in teams, or, even more important, when their results are to satisfy scientific standards and need to be replicable elsewhere, the analysts involved need to work not only together but also alike, or else their results will not be comparable. The coordination this requires is accomplished through the formulation of clear instructions for coders (see Chapter 7) to describe the same textual units in the same analytical terms, a data language (see Chapter 8). To ensure replicability, such instructions may include the following:

- A list of the qualifications that coders (observers, interpreters, judges) need for the task
- Descriptions of training procedures and instructional materials used to calibrate coders' conceptions
- Operational definitions of the recording and context units, and rules on how to distinguish them
- Operational definitions of the syntax (form) and semantics (meanings) of the data language (the categories or analytical terms) that coders are to apply in describing, translating, or categorizing each textual unit (Ideally, these definitions inform the cognitive operations that coders are asked to employ in reading and recording the texts. The definitions may be supplemented by examples of what should not be coded and why, including examples of what should not be included. See MacQueen et al., 1998; see also Krippendorff & Bock, 2009, Chapter 4.1.)
- Copies of the form(s) or electronic records to be used in creating records and entering data for processing: spreadsheets, examples of completed questionnaires, and initial tabulations

Typically, before these instructions are applied by several coders and to a large body of text, the analysts need to pretest them on a small sample of texts and then modify and retest them until they satisfy reasonable reliability standards (see Chapter 12, section 12.6.4).

There are several well-known strategies for developing suitable coding categories and recording instructions. Unfortunately, many content analysts use categories that are uniquely tailored to available texts, in effect starting each content analysis from scratch, almost in the spirit of text-driven approaches. Although this strategy eases the coding task and increases reliability, it creates content analyses whose results are not comparable with each other, and therefore rarely advance theory. Although ingenuity is always welcome, content analysts who rely on conceptualizations that have proven successful elsewhere have a better chance of contributing to existing knowledge.

A second strategy that many analysts use is to rely on the recording instructions of published content analyses with similar aims. I have mentioned various systems of categories in Chapters 7 and 8, and readers can find other examples in the works of authors such as Berelson (1952), Holsti (1969), Weber (1990), Gottschalk (1995),

Roberts (1997), Riffe, Lacy, and Fico (1998), and Neuendorf (2002), and in *The Content Analysis Reader* (Krippendorff & Bock, 2009). In addition, there are the categories built into computer programs, especially dictionary approaches (see Chapter 12, section 12.5.1). Some content analyses rely on only a few variables, whereas others define very many. Some require only a page of instructions; the instructions for others fill whole books (e.g., Dollard & Auld, 1959; Gottschalk & Gleser, 1969; Smith, 1992b). If existing coding categories have proven reliable and if they lead to answers to the research question under consideration, there is no need to invent a new scheme.

A third strategy is to draw from available literature on or theories of the context of the analysis. If the descriptive accounts or theories about this context can be operationalized into categories for coding texts, then analysts can gain immediate access to what the literature suggests the stable correlations are. This is the path that Osgood, Saporta, and Nunnally (1956) took repeatedly, for example, in developing evaluative assertion analysis. This analysis operationalized theories of cognitive balance, which led to such theoretical concepts as “attitude objects,” “connectors,” and “common meaning terms.” Where theories are unavailable, content analysis categories may be found in official classifications or taxonomies. If analysts need to describe occupational categories, it would make sense for them to consult official Federal Trade Commission listings or surveys of sociological studies of occupational status and prestige. For psychoanalytic research, the American Psychiatric Association’s *Diagnostic and Statistical Manual of Mental Disorders* (2000) is indispensable. As these categories are widely used, content analyses that use them can tap into empirical findings that are cast in these terms. Along the same lines, if a content analysis is to provide variables for testing particular hypotheses—for example, about participation in political deliberation—then the analysts might take their categories from previous research in this area, as Cappella, Price, and Nir (2002) did, relying for their argument repertoire index on Kuhn’s (1991) work on practical reasoning in everyday life. By deriving categories from established theories of the contexts of their analyses, researchers can avoid simplistic formulations and tap into a wealth of available conceptualizations.

In addition to having to be reliable, the categories of a data language should be tested, where possible, for their semantic validity (see Chapter 13, section 13.2.2). This is especially important for computer dictionaries, which, while perfectly reliable, may tag or transform text in incomprehensible ways. In Chapter 11 (section 11.4.1), I discuss several computer dictionaries and give one account of the development of categories suitable for computer processing (see also Krippendorff & Bock, 2009, Chapter 5.2). It may serve as a template for programming computer coding dictionaries.

#### 14.1.2.7 Selecting an Analytical Procedure

The best analytical procedures parallel what is going on in the context of the available texts. Figure 9.1 depicts the inferential component of content analysis as a procedural model of the presumed stable text-context relationships. Evaluative assertion analysis (see Chapter 9, section 9.2.3), for example, models the transfer of attitudes from common meaning terms to objects and from one attitude object to another. It operationalizes a set of psychological propositions, amounting to a

rather particular analytical construct. This analysis is appropriate only where attitudes are the target of research questions and the context conforms to how the process is theorized. Semantic network analysis (see Chapter 11, section 11.4.3) has a very different structure. It is an outgrowth of computational theories of cognition and is appropriate where these theories are proven valid.

The catalog of well-formulated analytical procedures from which analysts can choose is not very large. In making an informed choice from among several canned computer programs, assembling an analytical procedure from available components, or constructing one from scratch, content analysts are advised to consider three things: (a) be clear about the network of stable correlations in the analytical context chosen for the texts being analyzed, and (b) find out exactly how texts are regarded, processed, or transformed in the analytical procedures available, in order to (c) select the procedure whose computations or analytical steps provide the best model of the network of stable correlations and are therefore most likely to yield valid answers to the research questions.

Selecting among analytical procedures is not easy. Analysts should be aware that promoters of text analysis software tend to overstate their claims about what their software does—for instance, promising that it can extract concepts from text when all it does is calculate statistically outstanding word co-occurrences. For example, on close examination, claims that a software package mines content, models text, or develops and tests theories automatically often boil down to disappointingly simple procedural tricks that are far removed from what these claims suggest. Similarly, theorists often generalize the analytical powers of their own projects beyond available evidence. And, what is even more disheartening, most canned computer programs seal the assumptions built into their algorithms against outside inspection. These are some of the difficulties that analysts face in trying to make informed choices.

The LIWC (Linguistic Inquiry and Word Count) software (Pennebaker, Francis, and Booth, 2001; Pennebaker & Stone, 2001) is a notable exception. It is a dictionary approach (see Chapter 11, section 11.4.1). It clearly documents its rich, multi-dimensional dictionary, offers users numerous choices, and allows them to apply their own analytical constructs on top of the frequencies of dictionary entries it provides (for an application, see Krippendorff & Bock, 2009, Chapter 7.7).

#### 14.1.2.8 Adopting Standards

Given that the answers to content analysis research questions are inferences from texts about not-yet-observed phenomena, these answers are always of hypothetical validity. Standards serve to limit the uncertainty associated with such answers. This uncertainty is a function of three elements:

- The nature of the context of the texts being analyzed
- The extent of the analysts' knowledge of the text-context correlations
- The care with which the analysis is conducted

The true nature of the context from which texts are taken is not really under the analysts' control. Some contexts are highly institutionalized, whereas others are

open and situationally determined, if not deliberately ambiguous. In some, the connections between texts and the answers to research questions are linear and direct; in others, these connections are chaotic or probabilistic. This limits the certainty of the inferences that analysts can make. There is a danger that content analysts may oversimplify the construction of the context they end up utilizing in their analysis.

Knowledge of these correlations is another matter. No content analysis can be justified without some knowledge of this kind. But complete ignorance rarely exists. Content analysts are competent readers, at least in their own language, and do not let pass what seems incomprehensible to them. Beyond the ever-present face and social validity, content analysts may test for sampling, semantic, structural, and functional validity (see Chapter 13) and argue for the validity of their findings from the strengths of these tests, weaving information from appropriate literature and their own practical experiences into their rhetoric. In arguing for the validity of content analysis results, both proponents and critics rely on scientific standards of plausible reasoning. Such standards, although permeating discussions of scientific accomplishments, may not be quantifiable, attesting to the rhetorical nature of scientific research.

The third source of uncertainty, carelessness in conducting an analysis, shows up in at least two ways: as unreliability at the front end of an analysis, where it is measured with the help of suitable reliability coefficients (see Chapter 12), and in the way an analysis is designed to proceed to its result, also called internal validity (see Chapter 11). The recording/coding phase of content analysis is especially vulnerable to disagreements among coders, which show up in reliability tests.

But how high should standards be set? Obviously, when the results of a content analysis affect the life or death of a defendant in a court of law, when major business decisions with large price tags are based on them, for example, or when whole populations, say during wartime, are affected, standards need to be significantly higher than for scholarly work where the most that is at stake is the content analyst's reputation (see Chapter 12). The attainment of higher standards, although always desirable, tends to be more costly. It may require more careful preparatory investigations (see Chapter 4, section 4.2), a larger body of data (see Chapter 6), more sophisticated techniques of analysis, and so on. Content analysts may not wish to undertake projects whose standards are beyond their reach, both in terms of the needed resources and in terms of their responsibilities for disastrous consequences of mistaken inferences from texts.

Standards for sampling, semantic, structural, and functional validity should be related to the level of validity demanded of the results. To decide on such standards, researchers may want to work backward from how certain, general, or selective the results need to be to how often-unavoidable imperfections can affect them.

The relationship between reliability—a function of the agreement between two or more analytical processes, coders, or devices (see Chapter 12)—and validity (see Chapter 13) is quite transparent. High reliability is a prerequisite of high validity but does not guarantee it. Computer content analyses, while perfectly replicable, can be completely invalid. In Chapter 12 (see section 12.6.4), I discuss reliability standards and a way to test the consequences of unreliable data. The evaluation of the design of a content analysis, internal validity, is more qualitative in nature.



### 14.1.2.9 Allocating Resources

Content analysts have much to organize: analytical procedures, personnel, and scarce resources. Some activities may be reserved for the principal investigator, whereas others may be delegated to assistants, requiring training and instructions, or to professional research companies. Some must be executed sequentially—for example, the sampling of texts will have to take place before their coding, and coding must be done before analysis—and others may be done in parallel. Some take up short moments of time (e.g., running a computer program); others may be tedious (e.g., the reading and manual coding of text, most preparatory work, and the cleaning of dirty data). Limited resources—whether in qualified personnel, analytical devices, or funds—can impose organizational constraints on a project as well. Unless a content analysis is small and exploratory, analysts have to develop ways to organize their work efficiently.

There are numerous tools available to help analysts organize the processes of research. Most of these tools analyze the interconnected activities in a research project as a network. In such a network, arrows represent activities that one person or group can perform. By associating times and costs with each arrow, researchers can calculate needed resources, see potential bottlenecks, assign people to parallel or sequential activities, and estimate minimum and maximum amounts of time to completion.

Among the planning tools that content analysts may find useful are flowcharts such as those used in computer programming, the Program Evaluation and Review Technique (PERT), the Critical Path Method (CPM), and Gantt charts (interested readers can find information on all of these on the internet). These methods enable researchers to find the least expensive or the fastest paths to achieving research results and match available skills and resources with possible ways of organizing the analytical work.

### 14.1.3 Method-Driven Analyses

Method-driven analyses are suspect when they are motivated by what Abraham Kaplan (1964, p. 28) calls the “Law of the Instrument”: When a child discovers how to use a hammer, everything seems to be in need of hammering. Analogously, when researchers get hooked on one analytical technique, when they become experts in its use, they may well end up applying that technique to everything in sight—and not without pleasure. Technologies have this attraction, and content analysts, especially those employing computer-aided text analysis software, are not immune to this lure. Typically, mastering any reasonably complex analytical technique requires analysts to invest so much of their time that they find it increasingly difficult to shift gears, even to see alternatives outside their expertise. Instead of starting from real-life problems, content analysts can be tempted by this technological expertise to look for areas of research where their preferred methods are arguably applicable. This raises the possibility that the insights gained from method-driven analyses are more reflective of what particular methods can produce than of how the objects of inquiry operate.

On the positive side, when researchers conduct method-driven content analysis, especially with validity concerns in mind, they simultaneously expand their

method's areas of application while encountering its limitations. For example, the use of CatPac, a software package that fuses two theories—the self-organization of neuronal networks and the movement of diverse social objects within abstract spaces (Woelfel & Fink, 1980)—migrated from tracking advertising and public relations campaigns to optimizing mass-media messages to analyzing communication, bibliographic, and vocabulary networks in social settings (Barnett & Doerfel, 1997). CatPac is now known largely as a clustering program for qualitative data, for texts in particular. In the path it took to arrive at this point, the software encountered failures and critics but also demonstrated successes and gained proponents. It found its niche.

Method-driven content analyses face fewer design issues than do problem-driven analyses, largely because once a method is chosen, analytical options are limited. For example, in CatPac, recording units are unalterably defined as character strings, usually single words. Instead of ascertaining their meanings, CatPac applies an algorithm directly to these words. It clusters words using information about their co-occurrences within specified stretches of text. The clusters resulting from such an analysis are interpreted as representing conceptions in the minds of speakers, in the culture of an organization, or in the public at large. Proponents of CatPac consider this automatic mapping to be the software's most important virtue, whereas critics miss the use of human intelligence. Once a method is chosen, the research questions that can be answered are usually fixed. In CatPac, they concern the forming of clusters of textual units, interpreted as processes of conceptualization on several levels of abstraction, and this process is thought to be manifest in the sources of the analyzed texts.

In the design of method-driven content analyses, usually only five preparatory steps remain for analysts to accomplish:

- Locating and sampling relevant texts
- Ascertaining stable correlations in the contexts of these texts
- Preparing the texts in method-specific and context-sensitive ways
- Adopting standards or criteria
- Allocating resources

Locating and sampling relevant texts in method-driven analyses is less an issue of finding texts that correlate with the answers to a research question than one of locating texts that are easily amenable to processing by the chosen method.

Regarding the second step, method-driven content analysts are less inclined to explore correlations in a context for the directions an analysis could be taking than they are to ascertain whether the correlations surrounding the texts are compatible with the assumptions built into the chosen method. To continue the example above, CatPac assumes one-to-one relationships between single words or phrases and concepts in the minds of speakers. CatPac users are advised to examine literature or other sources of knowledge to determine whether its assumptions are warranted, whether the social phenomena of interest warrant this one-to-one relationship, and, more important, whether co-occurrences in texts are valid determinants of their meanings.

The preparation of texts in method-driven analyses resembles the previously described development of recording instruments, but not necessarily for use by coders. In CatPac, for example, analysts prepare a text by removing words deemed irrelevant, largely function words, and eliminating mere grammatical variations through stemming or lemmatization. Other text analysis software packages distinguish between go-words and stop-words, apply dictionaries that assign tags to words or phrases by which they are subsequently recognized, or parse sentences into components (see Chapter 11). Less automated analytical techniques, such as contingency analysis, require extensive manual editing of text. Justifications for these text transformations rely heavily on the analysts' judgments of what is relevant and what is not. Semantic validity is one applicable standard, reliability for manual editing is another, and computability by the method is a final and, in this approach, often primary criterion.

The adoption of standards in method-driven analyses essentially follows the arguments presented above, although some may not apply. For example, in computer analyses, reliability is not an issue. However, because computers are unable to comprehend texts the way humans do, semantic, structural, and functional validities are all the more important. Many computer aids fail semantic validity criteria but make up for this defect by processing huge volumes of text.

Finally, researchers must allocate their resources whether they are conducting problem- or method-driven analyses. The use of computational methods, once mastered, tends to be much cheaper than the use of methods that are chosen for their ability to answer method-independent research questions.

Method-driven content analysts tend to justify their use of methods by vindication (see Chapter 13). A method of analysis is vindicated when it consistently produces interpretable results. When a method crosses the limits of its usefulness, it produces obscure and uninterpretable results. The users of CatPac have had such experiences. Unfortunately, many researchers are hesitant to report their failures, even though content analysts can learn more about the limits of particular analytical methods from their failures than they can from their successes.

## **14.2 Writing a Research Proposal**

---

A research proposal puts forth the plan of a content analysis for consideration by a sponsor, dissertation committee, or teacher—someone who is able to grant permission, provide resources, or command time for the content analyst to engage in the proposed inquiry. As such, a proposal has both a rhetorical function and a contractual function.

### **14.2.1 Rhetorical Function**

The rhetorical function of the research proposal is to convince the sponsor(s) of two things:

- That the proposed research is worthwhile or beneficial to the sponsor
- That proponents of the research project are capable of delivering what they propose

In academic research, scholars tend to accomplish the first of these by citing relevant literature to demonstrate a gap in knowledge or in method that the proposed research can be expected to narrow. Ideally, this gap is of social significance, widespread, and instrumental to other advances, not just of personal interest to the researchers. However, all funding agencies have their own missions, which are manifest in their histories of funding certain research projects and not others, just as scholars in positions to approve research proposals have their theoretical concerns and epistemological commitments. For a proposal to succeed, it needs to address these. In applied research, clients tend to seek information with practical implications. To be successful, proposals should demonstrate that the benefits of the research outweigh its costs.

Often, researchers face competing expectations when they set out to write a research proposal. Such expectations may take the form of conflicting criteria from different departments of a funding agency, differing perspectives on the part of the members of a dissertation committee, or hidden agendas being pursued by decision makers who cannot talk about them publicly. In such a case, the researchers' best strategy is to write a proposal that enrolls all decision makers into the project, giving each a reason to support it. In dissertation research, this sometimes means that a student needs to write one chapter for each committee member. A commercial research proposal may want to show how each stakeholder in the proposed research could benefit from its results, or at least not be disadvantaged by them, and perhaps even how all stakeholders could be brought together on account of the proposed research project.

Past accomplishments are clearly the best recommendations of researchers' abilities. Sponsors examine proposals for researchers' academic degrees and lists of their publications as well as reviews of the analysts' previous research, especially those written by reputable critics, and letters of support from respected authorities. Researchers without relevant previous research records may compensate for this deficiency by providing compelling literature reviews in which they demonstrate familiarity with both the issues involved and how other researchers have solved or failed to solve similar research problems.

A research proposal does not merely discuss issues, however. It also needs to spell out the steps that the researchers intend to take and explain why. Indeed, probably the most convincing demonstration of the researchers' ability is a detailed research plan that evaluators can follow and critically examine for its likely success. The proposal should also report on any preparatory work the researchers have completed, spelling out or suggesting how the challenging problems of the proposed research are being solved. In content analysis, this often means that researchers should present evidence of the reliability of the proposed recording instructions and the capability of the chosen analytical or computational techniques to advance the aim of the research, as well as accounts of the relevance to the research question of the texts to be sampled.

### 14.2.2 Contractual Function

A research proposal, once approved, entails the expectation that the sponsor or funding agency will provide what it has agreed to make available—financial

resources, organizational help, or legal support—and that the researchers will deliver what they have proposed. The approval of a proposal creates contractual obligations whether the proposed research is intended to qualify an individual for an academic degree, to contribute theoretical insights, or to provide business or political intelligence.

One of the characteristics of scientific research is that its results cannot be guaranteed before the research is completed. Just as researchers who propose to test a hypothesis must consider evidence in favor of that hypothesis as well as against it, so must content analysts who propose to answer a certain research question keep its possible answers open for the analysis to decide among them. Nature does what it does, and texts do not always yield what sponsors like to see and analysts hope to show.

In fact, lack of “cooperation” of texts often stimulates new insights and opens unanticipated turns, which brings us to the second peculiarity of scientific research: serendipity. A research proposal must outline at least one path to answering the research questions but at the same time preserve the analysts’ ability to deviate from that path when unanticipated shortcuts become apparent, when new methods turn up, or when unforeseen findings shift the focus of attention to something even more exciting—provided the research objective stays within expectations and scientific standards are not compromised.

The inability to guarantee particular results and serendipity in the conduct of research can be anathema to funding agencies that have vested interests in preferred outcomes. Proposal writers may need to address this peculiarity of research and convince the sponsors that all legitimate research questions have alternative answers, or formulate their research questions so that the sponsors see virtue in the possibility of less preferred answers.

### 14.2.3 Outline for a Research Proposal

A typical proposal for a content analysis includes all of the following parts:

- *A statement of the general epistemic or methodological issue* that the proposed analysis will address: what that issue is and why and to whom it is significant
- *A review of available literature on the context* in which this issue resides, showing the kinds of questions that have been asked and answered, the kinds of research methods previously applied, and what has worked and what has not, including the analysts’ own research or experiences, if relevant
- *A formulation of the specific research questions* to be answered by the proposed research, which should be embedded in an account of the analytical framework adopted, the *context chosen by the analysts* to make sense of these questions, and the *body of text* by means of which the analysts expect to answer these questions (see Chapter 2)
- *A step-by-step description of the procedure to be followed*, including accounts of any preparatory research already undertaken or to be carried

out, the hypotheses to be tested (see Chapter 4, section 4.2) as well as how and why they are to be tested, the proposed analytical steps (Figure 4.2), and the standards adopted for each. This description should cover the following:

- The *units of analysis* (see Chapter 5) proposed, defined, and distinguished, and what they respectively contain and omit from the body of text
  - The *sampling strategies* (see Chapter 6) to be used, where the population of relevant texts is located, how easily available it is, criteria for adequate sample sizes, methods for correcting self-sampling biases, and the sampling validity to be achieved (see Chapter 13, section 13.2.1)
  - The *recording/coding categories* and *data language* (see Chapters 7 and 8) to be used, whether in the form of instructions to coders, available computer dictionaries, or queries for text mining (see Chapter 11) or to be derived from theories, literature (see Chapter 9, section 9.3), or the texts themselves (see Chapter 1, section 1.7; see also section 14.1.1); what these categories preserve or omit; semantic validity to be achieved (see Chapter 13, section 13.2.2); the reliability to be guaranteed (see Chapter 12); and the results of any pretests of the recording instructions
  - The computational (statistical or algebraic) *techniques for reducing* or summarizing the body of recorded text and the justifications for these techniques relative to what is known about the context of the texts
  - The *inferential procedures* that the analysts will ultimately use to answer research questions from texts (see Chapter 9): the analytical constructs that underlie these and any evidence for their structural and functional validity (see Chapter 13, sections 13.2.3 and 13.2.4), available computer programs to be used, and evidence of their previously established correlative or predictive validities, if any (sections 13.2.5 and 13.2.6)
  - How and to whom the research results are to be made available: the *narrative forms* in which the answers to the research questions will be presented, using numerical arrays, graphic demonstrations, or computed indices, for example, and planned publications, presentations to conferences, or reports to sponsors; the *kinds of conclusions* drawn from the results, whether they are expected to advance theoretical knowledge, make recommendations for actions, or settle an issue; and a *critical assessment* of the uncertainties that are likely to remain associated with the larger issues that the results are to address
- *An account of the specific time periods and resources needed* to complete the proposed analysis (the costs of personnel, equipment, and outside services) presented in the form of a timeline of the phases of research showing the milestones to be achieved and the resources needed at each phase
  - *A list of references to cited literature* that conforms to whatever style manual the sponsor of the proposal accepts and includes entries only for available publications

- *Appendixes* containing material pertinent to the proposal but not central to the potential sponsor's understanding of what is being proposed, such as examples of the kinds of texts to be analyzed, lists of texts to be sampled, the proposed categories of the analysis, its data language and/or recording instructions if already available, preliminary reliabilities achieved so far, specifications of the software to be utilized, preliminary analytical results, and testimony by experts supporting the proposed research

## 14.3 Applying the Research Design

---

Ideally, the work involved in carrying out a well-designed content analysis becomes routine—even when a content analysis is qualitative, text-driven, or hermeneutic and exploratory to begin with. With all intellectual and methodological problems solved during the design phase, the analysis could be turned over to a research organization. In practice, however, problems are bound to emerge: Needed texts may turn out to be unavailable, scarcely relevant, or biased by incorrigible self-sampling practices of their sources, and software may turn out not to work as expected. However, the most frequent disruptions stem from the inability to meet accepted standards of reliability (see Chapter 12). As solutions to such emerging problems cannot be specified in advance, short of discontinuing a content analysis altogether, researchers may have to be prepared to go back and modify the defective parts of a research design (see Figure 4.2), keeping the overall research objective in mind. It is not unusual for a content analysis to need several iterations of locating unreliabilities, correcting their causes, and repeating these steps until applicable standards are satisfied. Often researchers get stuck with coding instructions that turn out to be unreliable, and they introduce small modifications instead of radically reconceptualizing their approach. In addition to unreliabilities, new empirical findings and literature about the context of the analyzed texts can prompt reconstructions of the context an analysis was presupposing. Content analysts cannot be afraid to go back to solid ground, as the issue always is to select potentially valid answers to the research question.

## 14.4 Narrating the Results

---

Research proposals are written to convince sponsors, but research reports typically are addressed to other readers. Also, research reports usually go far beyond the mere statement of findings of fact. Such reports account for how analysts have accomplished what they set out to do; describe to which literature, judgments, or decisions the research results contribute; and raise questions for further exploration. In highly institutionalized settings—such as laboratories or public opinion polling organizations—where research questions are codified, researchers are well-known, and analytical procedures are well established, research reports may be limited to information about where the analyses deviated from the typical. Generally, a

research report should offer details sufficient to convince at least three kinds of addressees of the importance of the results:

- The sponsor, agency, or client who approved and/or supported the research
- The content analysts' peers in the scientific community
- The public at large

These addressees may have conflicting agendas, which may have to be attended to separately.

Sponsors are interested, first, in whether the researchers fulfilled their contractual obligations. A research report needs to demonstrate that they did, and, where applicable, should also justify where and why they deviated from the original proposal. Second, funding agencies are keenly aware of the publicity they gain from having supported worthwhile research and often look for the social validity or political significance of the results. Political or commercial clients might be more interested in the benefits they can reap from the research, whereas academics may look for the knowledge advanced by the results. A research report may need to address these concerns. Third is the issue of empirical validity: Can the analysts' claimed results be trusted? Nonscientific users may be inclined to accept content analysis results on account of the analysts' scientific credentials or the reputations of the institutions with which they are associated. Users in the legal and scientific communities, by contrast, may approach a research report more critically, looking for and needing to find relevant details and evidence in support of findings.

For members of the scientific community, the ability to reproduce research results elsewhere is the most widely used standard. As the burden of proof is laid on the researchers, the research report must provide convincing evidence that the results can be reproduced. Measures of the reliability of the potentially unreliable components of an analysis—recording/coding, for example—may provide the needed assurances. Many scholarly journals require evidence for the reliability and statistical adequacy of research findings. But reliabilities merely put an upper limit on the potential validity of research results, as discussed above.

Because the answers that content analyses yield are obtained through abductive inferences, analysts need to establish the validity of their inferences. They may accomplish this by retracing the analytical steps taken and justifying each step in terms of whether it models or represents what is known about the context of the texts. Analysts can use sampling, semantic, structural, and functional validities (see Chapter 13) to support such arguments, especially when their results seem counterintuitive. However, content analysts should expect that even plausible research results may be received with a healthy dose of suspicion. The face validity they perceive may not be clear to all those who are touched by the research findings, forcing the analysts to go the extra mile to provide arguments that are compelling, even for their critics, even if the validity of their findings may seem obvious within their community of peers.

When research results are published, they enter the conversations of diverse readers, if not public debate. To succeed in the view of the public, analysts may want to provide compelling narratives that emphasize the impacts their findings will



have on the readers' lives, using concepts, comparisons, or metaphors that are meaningful in the readers' own worlds and not misleading.

### 14.4.1 Outline for a Research Report

Generally, a research report should contain the following parts (many of which are also found in the research proposal; see section 14.2.3):

- A *summary or abstract* of the research for decision makers who have little time to concern themselves with details (This part of the report often decides whether readers will continue with the remainder.)
- A table of contents
- A statement of the epistemic or methodological issues that informed the analysis
- A *review of the literature* on the context in which these issues reside
- An *account of the framework adopted* for the content analysis, including the research questions addressed (with an explanation of the sponsor's or analysts' interest in these questions), the texts analyzed to answer the questions, and the context chosen to justify the analysis (see Chapter 2)
- A *description of the research design* actually followed, including the preparatory research undertaken (see Chapter 4), any complications encountered in the process, and how emerging problems were solved, with specific information in the following areas to enable critical evaluation of the process:
  - The *body of texts* sampled: what it consists of, what motivated the analysts to choose it, by which strategy it was selected (see Chapter 6), and how the analysts dealt with biases (see Chapter 13, section 13.2.1)
  - The *data language* used (see Chapter 8): the system of descriptive categories and measurements the analysts employed to bridge the gap between raw texts and the computational techniques applied
  - The *units of analysis* (see Chapter 5): their operational definitions, how they were used in the process, and what they preserved and ignored from the texts
  - The *recording/coding process*: whether built into computer dictionaries or search queries (see Chapter 11) or enacted by human coders (see Chapter 7), the reliability (see Chapter 12) and, where evaluated, the semantic validity (see Chapter 13, section 13.2.2) of each variable
  - The *computational* (statistical or algebraic) *techniques employed to summarize, simplify, or reduce* the volume of records obtained from the body of texts
  - The *inferential techniques* (computer programs [see Chapter 11] or other analytical procedures [see Chapter 9]) utilized to answer the research questions and, where available, evidence of structural validity (see Chapter 13, section 13.2.3) or functional validity (section 13.2.4)
  - The *research results (the answers to the research questions)*: in the form of data files, summaries, statistical accounts, propositions of a factual nature, recommendations for actions, or judgments (all in view of their potential validity), crafted in a compelling narrative that the anticipated readers of the report can easily understand

- A *self-critical appraisal*: of the analysis (Did it really yield something new?), of the time and resources spent (Was it worth the effort?), of the methods used in relation to other methods (Could there have been more appropriate techniques?), of the computational procedures used (Did they accomplish what was expected of them?), and of the meta-accomplishments of the analysis (Did it contribute to the content analyst's field? For example, did it raise new questions for further explorations?)
- *Additional matter*
  - The complete list of references to cited literature
  - Appendixes containing materials that enable interested readers to read beyond the report, such as the recording instructions, computer dictionaries used, the list of reliabilities obtained, and tables of numerical findings, even data that could be used by other researchers (Scientific research is open, and data may need to be made available for reanalysis elsewhere.)
  - Acknowledgments of contributors to the research effort (All research projects proceed in networks of interpersonal relations. Because assistants, coders, consultants, advisers, librarians, and teachers do not expect to be acknowledged in official research reports, it is a gesture of generosity when researchers name those who have been involved.)

## CHAPTER 10

# Analytical/Representational Techniques

*Methods in content analysis largely address the making and processing of data and the application of analytical constructs that preserve some of the data's meanings, leading to valid inferences. This chapter discusses ways in which researchers can represent the results of content analyses such that they may recognize patterns and discover new ways of exploring their findings. Such representations are informative relative to often-implicit standards, several of which are reviewed in this chapter.*

**A**fter texts have been recorded and analytical constructs have been applied, the content analyst needs to do the following:

- Summarize the inferences from text so that they are easily understood, interpreted, or related to intended decisions
- Discover patterns and relationships within findings that an unaided observer would otherwise easily overlook, to test hypotheses concerning various relationships
- Compare the findings with data obtained by other means or from other situations to support conclusions drawn from other research (multiple operationalism), to gain confidence in the validity of the content analysis at hand, to add another dimension to the intended inferences, or to provide missing information

In practice, these three tasks are not entirely distinct. They are not entirely unique to content analysis either. Much scholarly work, especially in statistics, is concerned with summarizing large bodies of data, making various comparisons, and testing statistical hypotheses. I cannot possibly review all techniques that content analysts might use, so I focus in this chapter on a few that benefit content analysts especially.

Moreover, I will not attempt to discuss these techniques in such detail that readers can replicate them—some require expertise found in common textbooks on research methods, and others are built into readily available statistical packages. Rather, my aim in this chapter is to suggest ways of analyzing and representing results tied to texts.

## Counts 10.1

---

Owing to the large amount of text that content analysts typically consider, counting parts of it is by far the most common technique used to reduce its volume to something manageable and presumably still comprehensible. The qualification “presumably” is intended to warn novices that counting always eliminates all relationships between the units counted and the simplicity of numerical accounts is always purchased at the cost of omitting the contextually complex meanings between the units in text. It is important not to conflate counting with the use of scientific methods. The early content analyses described by Berelson (1952) and justified by Lasswell (1965b) did just this. In this section, counting is taken as a practical way of coping with large volumes of text, as a method of data reduction, described in Chapter 4, section 4.1.1. With this use in mind, counting is justifiable only when the resulting frequencies mean something, or have something to do with the context of texts, including selecting quantitative answers to research questions.

Referring to the distinction between open and closed variables (introduced in Chapter 8, section 8.3), it is important to distinguish counts of mentions—that is, of textual units, words and phrases, that occur as such in a body of text, unanticipatedly—from counts of coded, categorized, scaled, or tabulated textual elements, whose terms are exhaustively defined by the analyst.

For content analysts, counting is justified only when the resulting frequency accounts can somehow be related to what a body of text means in the chosen context, whether it leads to answering a research question. Counting mentions is most elementary, and, when it comes to the use of computers, it is the easiest account obtainable. Not only because synonyms, grammatical variations of words, and typos are not distinguished in counts of mentions, such counts have difficulties keeping intact stereotypical phrases such as “content analysis,” “computer literacy,” or “museum of modern art” and can also include words with purely grammatical meanings. So-called word clouds are popular visualizations of frequencies of mentions. In 2005, Amazon.com added such accounts to its “inside the book” feature. They can be obtained from any text and for all languages with clearly distinguished words.

The following example, a word cloud of Chapter 1 of this book, contains the 50 most frequent words in the chapter (omitting function words) and depicts their frequency by the size of the type. As Mary Bock observes, “These counts make no scientific claims or inferences; instead, they offer their readers the possibility of unconstrained interpretations, based solely on the assumption that word frequencies mean something” (Krippendorff & Bock, 2009, p. 38). In the absence of inferences

analyses **analysis** analysts analyzing approach became berelson  
 communications computer **content** data development different  
 documents emerged example group interact interest MASS matter meanings media messages **newspaper**  
 numerous **political** press processing propaganda psychology public published qualitative  
**quantitative research** results **social** software study subject survey symbols **text** theory  
 used volumes War work world

---

## Word Cloud of Chapter 1

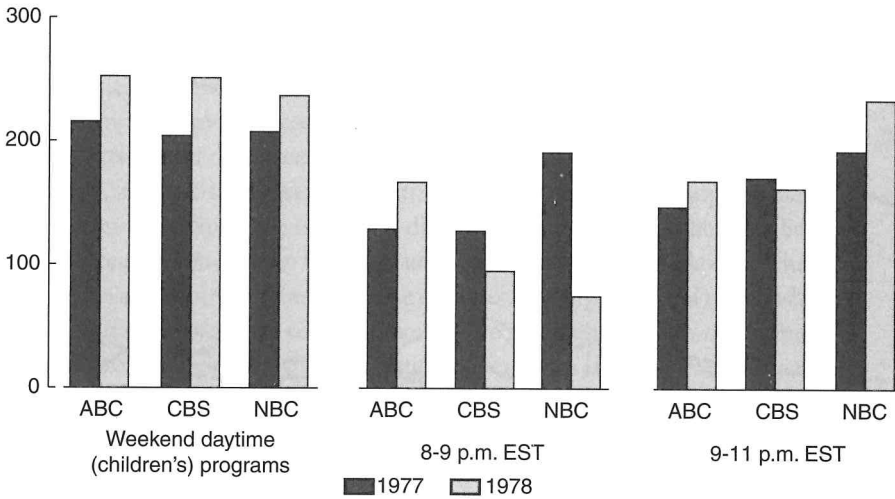
Source: Generated at <http://tagcrowd.com>; accessed June 20, 2011.

from text, Bock calls such accounts “impressionistic content analyses.” Those who have read Chapter 1 may find this visualization plausible, but from its word cloud, it would be impossible to reconstruct the content of this Chapter 1. Accounts of mentions amount to the lowest-level abstraction from text.

Counts of textual units in terms of coding categories of closed variables are more interesting. In fact, coding, categorizing, or scaling relevant textual recording units resembles acts of abstracting textual matter on which qualitative content analysts rely heavily. Although frequencies are often celebrated for their precision and simplicity, they should not be granted any special scientific significance. In comparing the results of counts in terms of coding categories, analysts typically apply several interpretive standards, often without being explicit about them. Content analysts should recognize and note their standards explicitly. I discuss three standards in the current section and one in the next.

Content analysts apply the standard of a *uniform distribution* when reporting that the frequency in one category is larger or smaller than the average frequency for all categories. The idea of bias in reporting, such as attending to one candidate for political office more than to another, exemplifies the implicit use of this standard. If favorable and unfavorable accounts were the same for both candidates, analysts would not call this bias and probably would not bother writing about it—except perhaps in surprise, because equality in coverage rarely happens. Figure 10.1, which is taken from Gerbner, Gross, Signorielli, Morgan, and Jackson-Beeck’s (1979) work on television violence, invites questions about such issues as why weekend children’s programs are so much more violent than other programs and which networks increased or decreased the violence in their programming from 1977 to 1978. When content analysts find observed frequencies noteworthy, then the frequencies deviate from what would not be noteworthy, and that usually means deviations from uniform distribution of categories. The bar graph in Figure 10.3, which comes from Freeman’s (2001) study of letters to auto industry shareholders, does not even show frequencies, displaying only deviations from their average, just what is significant.

When analysts observe changes in frequencies over time, they are likely to ask why some changes are irregular and deviate from what would be expected if



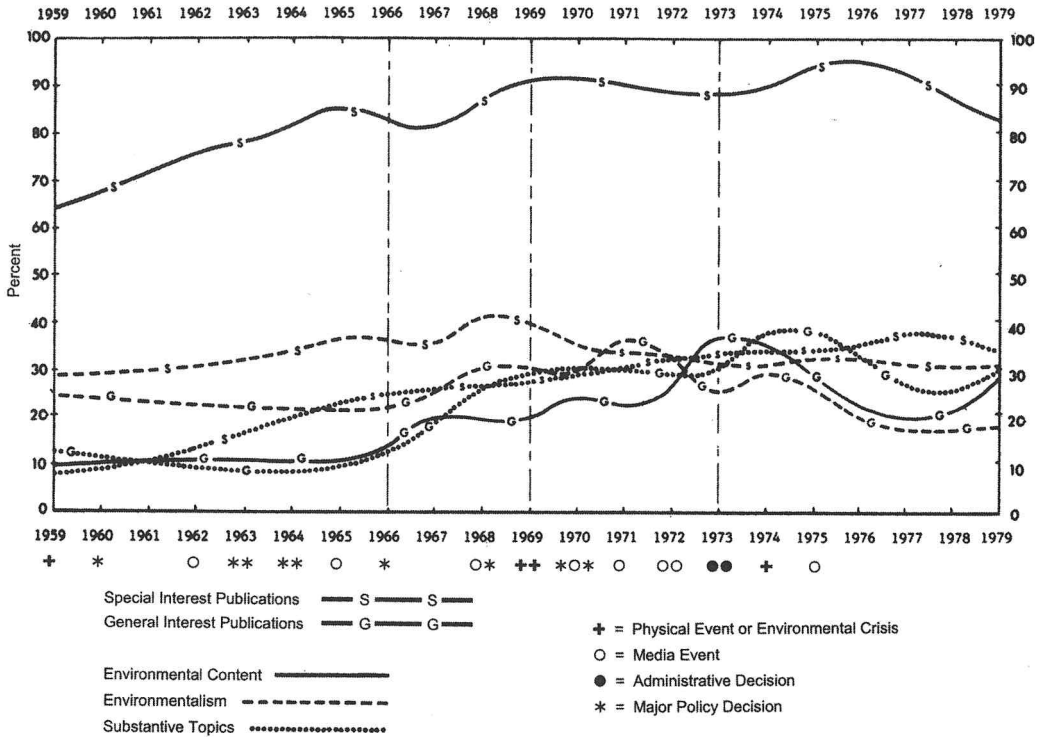
**Figure 10.1** Bar Graph Representation of Frequencies

Source: Gerbner et al. (1979).

changes were regular and predictable. They then refer to a *stable temporal pattern* as an interpretive standard; deviations from that pattern are noticed and considered important. Figure 10.2, which comes from an analysis conducted by Strodtzoff, Hawkins, and Schoenfeld (1985), shows trend lines for environmental content, environmentalism, and substantive content of special-interest and general-audience channels, largely magazines. The researchers also list four kinds of events in hopes that these might explain the deviations from the otherwise smooth increase over time.

Equally important and perhaps more typical in the content analysis literature is the standard of *accurate representation*, which is implied when an analyst notes that the relative frequencies differ from what would be expected if the data were a statistically correct representation of a population. This standard was introduced by Berelson and Salter (1946), who compared the population of characters featured in magazine fiction with the known demographics of the U.S. population. They found that minorities and poor people were all but absent in magazine fiction, and that popular heroes were overrepresented. Many critics of the mass media have noted that the population of television characters is not representative of the U.S. population or of the members of the mass-media audience. In early television research, analysts demonstrated this lack of representativeness especially for ethnic groups but also for people in particular occupations, people with low socioeconomic status, women in positions of leadership, and elderly people, and these studies were used to infer social prejudices, economic interests, and technological biases.

Whether a population of audience members is the appropriate standard against which the population of television characters should be judged is debatable. Many



**Figure 10.2** Trend Lines for Environmental Content, Environmentalism, and Substantive Content of Special-Interest and General-Audience Channels Juxtaposed With “Real World” Events

Source: Strodthoff et al. (1985, fig. 4).

popular figures, from film stars to television commentators, exist only in the media, not in any unmediated population, and there are good reasons popular talents are more likely to be shown on the screen and attended to than are ordinary people doing ordinary things. Yet, the application of the standard of accurate representation can have political consequences. For example, content analysis research in the late 1950s demonstrated the systematic underrepresentation of African Americans on U.S. television, and these research findings contributed to the eventual achievement of at least some racial balance on U.S. television.

## 10.2 Cross-Tabulations, Associations, and Correlations

The standard of *chance* is probably most common in statistical accounts of content analysis findings. It arises from analysts’ efforts to cross-tabulate the frequencies of several variables and to observe the frequencies of co-occurrences of values or categories rather than of simple categories. For example, a content analysis of 2,430 acts

performed by television characters yielded the observed/expected frequencies of co-occurrences shown in Table 10.1 (Brouwer, Clark, Gerbner, & Krippendorff, 1969). Simple frequencies say nothing about relationships between content variables. Table 10.1, for example, shows that good characters are the origin of most acts, a total of 1,125, followed by 935 acts by bad characters and 370 by neutral ones. Out of the 1,125 acts by good characters, most (751) are unrelated to the law. Although these are large frequencies, and they are far from being uniformly distributed, the actual numbers say little about the relationship between the favorable-unfavorable evaluation of television characters and their association with the law. If one is interested in the statistical relationship between two variables, one must compare the observed frequencies of co-occurrences with those obtained by chance. In cross-tabulations, frequencies are at the level of chance when all columns and all rows are proportional to their respective margins, which means that the marginal frequencies explain the distribution of frequencies within the table. In Table 10.1, the frequencies obtainable by chance are depicted in italics, directly below the corresponding observed frequencies in bold typeface.

**Acts Initiated by Fictional Characters** who are:

	Good	Neutral	Bad	
Associated with Law Enforcement	<b>369</b> <i>194</i>	<b>27</b> <i>64</i>	<b>23</b> <i>161</i>	419
Unrelated to Law	<b>751</b> <i>710</i>	<b>328</b> <i>233</i>	<b>454</b> <i>590</i>	1,533
Criminals	<b>5</b> <i>221</i>	<b>15</b> <i>73</i>	<b>458</b> <i>184</i>	478
Totals of	1,125	370	935	2,430 Acts

**Table 10.1** Cross-Tabulation of Frequencies of Acts Engaged in by Characters in Fictional Television Programming

Source: Brouwer et al. (1969).

What is noteworthy in such a table are co-occurrences of categories whose observed frequencies deviate significantly from what would be expected when variables were independent and co-occurrences were chance events. In Table 10.1, the largest observed frequency of 751 is also nearly as expected, 710, and thus does not contribute to the significance of the relationship between the two variables. In fact, when one uses a  $\chi^2$  test of this significance, the cells that make the largest contribution to this relationship are the four corner cells, which indicate the extremes of good and bad and of upholding and breaking the law. The differences between the observed and the expected frequencies in these cells tested statistically significant, and thus can be interpreted as supporting the statistical hypothesis that the *good guys* are more likely acting on the side of the law, whereas the bad guys are acting in opposition to it. I say “statistical hypothesis” here because the table shows that there are exceptions, although significantly fewer than chance.



Cross-tabulations are not limited to two or three variables, but they are more easily visualized and interpreted when the number of variables is small. Multivariate techniques are available for testing complex structures within multidimensional data (Esbensen, 2010; Reynolds, 1977).

When variables are nominal (an unordered set of categories), we speak of associations, as shown above, but when they consist of numerically ordered values, we speak of correlations. The standard of chance is common to both, but the use of correlation coefficients adds another standard to that of chance, the standard of *linearity*. Correlation coefficients are zero if data are as expected by chance, and they are unity when all data fall on a straight line, a regression line (see also Chapter 12, section 12.7). Otherwise, correlation measures the degree to which data resemble a regression line as opposed to chance. Above-chance statistical relations—associations and correlations—may be of two kinds:

- Within the results of a content analysis, as in Table 10.1
- Between the results of a content analysis and data obtained independently, as in Figure 10.3

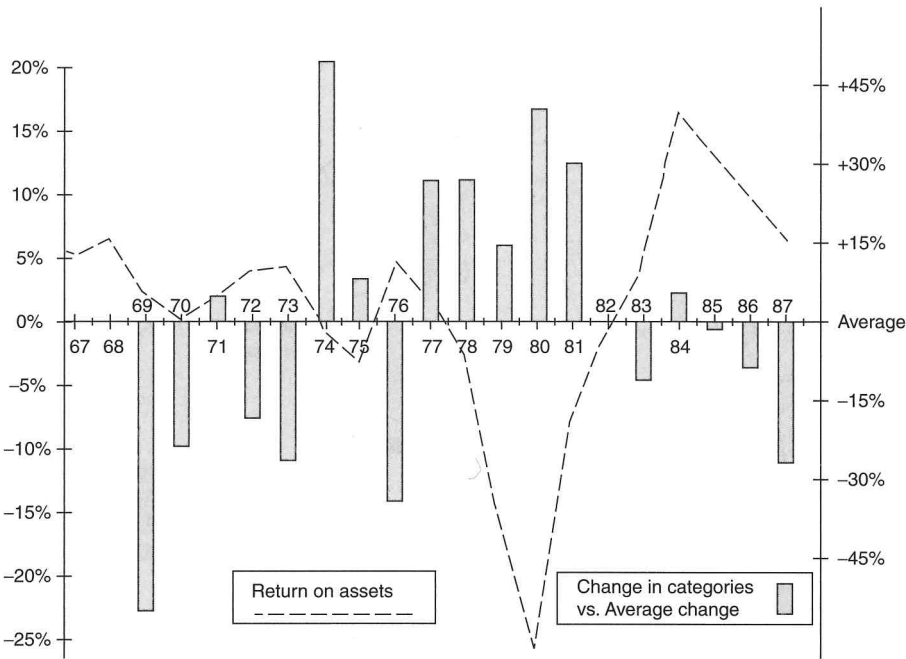
Because content analysts control the definitions of their variables, there is always the danger that the statistical relations within content analysis results are artifacts of the recording instrument. In Table 10.1, the positive association (good cops, bad criminals) is notable because the underlying relation could have gone in the other direction (bad cops, good criminals). But a positive association between, say, feminine-masculine personality traits (gender) and sex (its biological manifestation) is expected in our culture precisely because these two variables are semantically related. Association and correlation coefficients do not respond to semantic relationships between variables, and if such relationships do exist, these correlation measures are partly spurious and uninformative by themselves.

Correlations between the results of a content analysis and data obtained by other means are less likely so affected because the two kinds of variables differ in how the data are generated. Figure 10.3 comes from Freeman's (2001) study of U.S. auto industry letters to shareholders. Freeman compared the attention paid to a set of categories functional to a corporation in Chrysler's letters to shareholders with the company's return on assets and found a strong negative correlation between these variables.

## 10.3 Multivariate Techniques

---

The standard of chance underlies most multivariate techniques of data analysis and representation. Correlations are worth reporting only when the data deviate significantly from chance, ideally approximating linearity. One prominent technique is multiple regression analysis. It presupposes that the variables being analyzed are of two kinds: independent and dependent. The variation in the dependent variables is to be explained, and the variation in the independent variables serves as the explanation. Indeed, many questions that content analysts pursue are reducible to

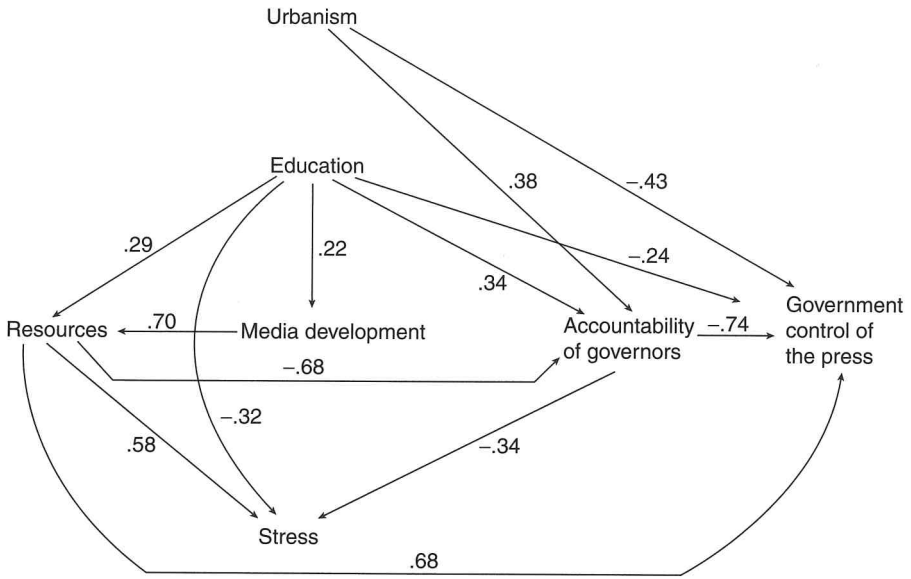


**Figure 10.3** Correlation Between Chrysler's Return on Assets and Year-to-Year Attention to Functional Categories in Chrysler's Letters to Shareholders

Source: Freeman (2001, fig. 5).

problems of regression. For example, which characteristics of novels predict their popularity? A clear answer to that question would please authors and publishers alike. Or which factors explain media content—government actions, interest groups, economics (advertising), technology, or artistic talent? Or which features of messages are effective in encouraging members of a target population to change their health care habits? The most common kind of regression analysis orders a number of independent variables according to how much they contribute to predicting the values of one chosen dependent variable.

Another multivariate technique entails the use of structural equations. Each variable is considered a dependent variable of all other variables. Only under certain conditions can such a network of multivariate correlations be interpreted in causal terms. Constraints of space prevent me from discussing these conditions here, but I must note that it is extremely difficult to establish causality from exclusively textual data. One important ingredient of the use of causal explanations is time. Figure 10.4 shows the results of a path analysis conducted by Weaver, Buddenbaum, and Fair (1985) that features regression correlation (beta) coefficients above .20 between variables whose relationship to the development of the media in Third World countries was suspected. Weaver et al. compared this path analysis with one that used the same variables but also included data concerning all countries and concluded that in



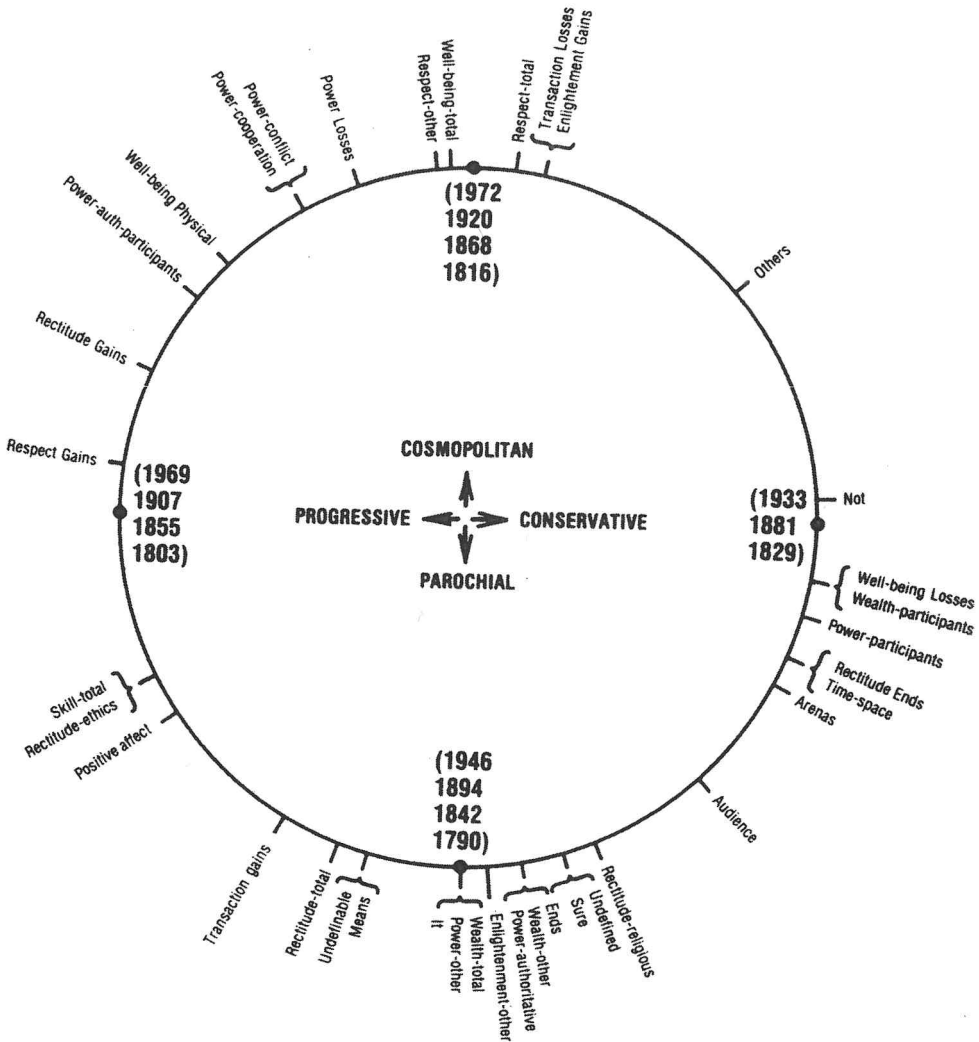
**Figure 10.4** Paths for Predicting Governmental Control of the Press for Third World Countries, 1950–1979

Source: Weaver et al. (1985, fig. 2).

most Third World countries, the media tend to be used to facilitate the functioning of the economy and to perpetuate the power of the rulers.

Correlational techniques are not limited to linear relationships, however. A good illustration of this is found in the work of Namenwirth (1973; Namenwirth & Weber, 1987), who analyzed the Lasswellian values (see Chapter 8, section 8.5.3) in speeches from the British throne between 1689 and 1972, covering the British mercantilist and capitalist periods. Over such a long period, fluctuations in the frequencies of values are to be expected, but instead of correlating these with external events that the British Empire had to face, Namenwirth considered values as expressing the workings of an autonomous cultural/political system in which the frequencies of one kind decline as others rise, in endless cycles. To test this hypothesis, he applied a kind of Fourier analysis to these fluctuations. A Fourier analysis decomposes the complex fluctuations of a measure—whether of waves of light or of economic activity—over time into a series of additive sinus curves. Namenwirth identified at least three concurring cycles that turned out to explain much of the variance in the data: a 146-year cycle, a 52-year cycle, and a 32-year cycle.

Figure 10.5 depicts the internal structure of the 52-year cycle. Categories of values that peak are listed at the rim of the circle. Accordingly, the “6 o’clock” position, corresponding to the years 1790, 1842, 1894, and 1946, witnesses concern about the poor performance of the economy, prevailing (un)certainties, and search for knowledge (enlightenment). In the “9 o’clock” position, rectitude and respect gain and concerns for social welfare and conflict grow. In the “12 o’clock” position, welfare



**Figure 10.5** A 52-Year Cycle of Values Found in Speeches from the British Throne, 1689–1972

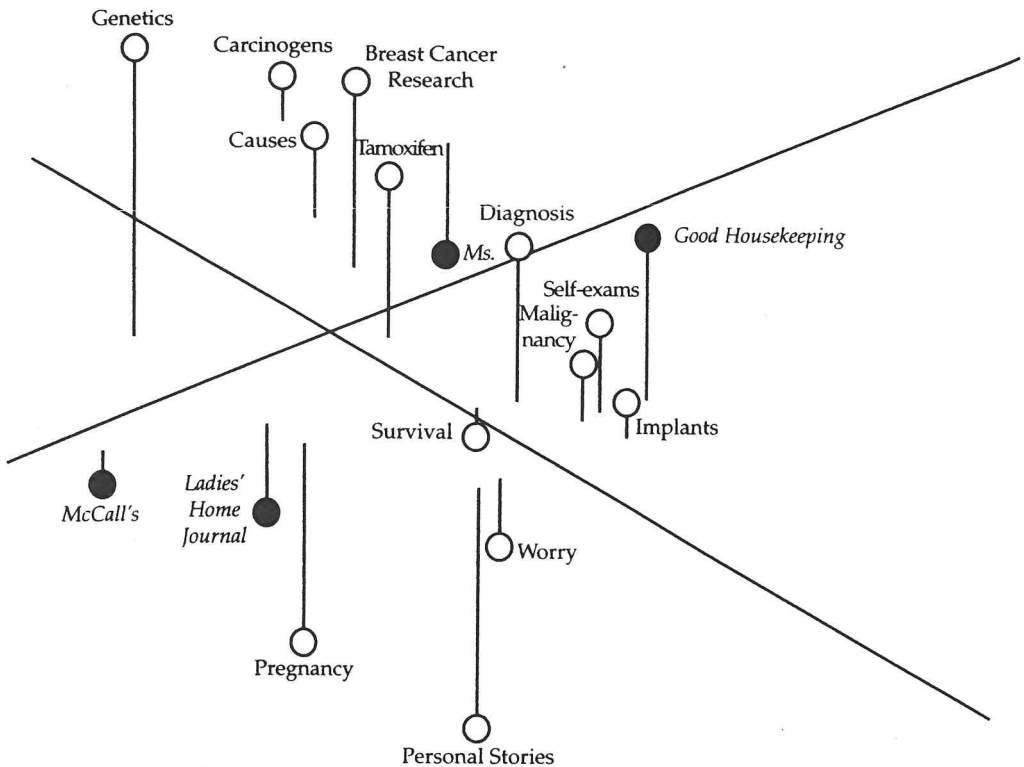
Source: Namenwirth and Weber (1987, p. 139, fig. 5.5; also in Namenwirth, 1973).

and respect reach their peaks, and enlightenment co-occurs with an international orientation. At the “3 o’clock” position, wealth, trade, and conflict become issues, and well-being is feared to degenerate. In the center of this figure, Namenwirth summarizes this dynamic as a thematic progression from parochial to progressive, to cosmopolitan, to conservative, and back to the beginning. This interpretation loosely follows Parsons and Bales’s (1953) theory suggesting that every society, facing four functional requirements, cycles from an expressive phase to an adaptive phase, then to an instrumental phase, then to an integrative phase, and then back to an expressive phase, and so on. One could describe this technique as one of curve fitting. Here the usual linearity assumptions of correlation coefficients are replaced by sinus curves.

## 10.4 Factor Analysis and Multidimensional Scaling

Factor analysis, a favorite method of behavioral scientists in the 1960s and 1970s, is a way to summarize the correlations among many variables by constructing a space with fewer dimensions in which these data might be represented with a minimum of loss in explanatory power. It computes a set of hypothetical and ideally orthogonal dimensions or variables and offers measures of how closely the original variables are correlated with these. These correlations (of the original variables with the virtual dimensions) provide clues that help analysts to make sense of the virtual dimensions. This is the path that Osgood (1974a, 1974b) took to obtain what he called the “basic dimensions” of affective meaning. He used data in the form of numerous semantic differential scales and found three basic dimensions that explain between 50% and 75% of the variance. After examining which of the original scales correlated highly with these, he called them the “evaluative” (good-bad), “activity” (active-passive), and “potency” (strong-weak) dimensions of affective meaning (see also Chapter 7, section 7.4.4).

Whereas factor analysis reduces the dimensionality of the original data while trying to preserve their variance, multidimensional scaling (MDS) reduces the



**Figure 10.6** Three-Dimensional Representation of the Frames Used by Four Women's Magazines to Discuss Breast Cancer

Source: Andsager and Powers (1999, p. 541, fig. 1).

dimensionality of the original (geometric) distances between data points, trying to preserve their positions relative to each other. It requires data on how far apart pairs of elements, concepts, and even variables are. The analyst can fulfill this condition in various ways, such as by measuring differences, dissimilarities, disagreements, dissociations, or lack of co-occurrences between all pairs, whether using objective measurements or subjective judgments. Even correlation coefficients can be and have been converted into distances and subjected to MDS techniques.

MDS starts out with a space of as many dimensions as there are data points, which usually escapes human comprehension. It then attempts to remove one dimension at a time, so as to represent the data in fewer dimensions with a minimum of adjustments to the distances between the data points—much as when one attempts to represent a three-dimensional distribution of points in two dimensions. Figure 10.6 displays the MDS results of a content analysis conducted by Andsager and Powers (1999), a three-dimensional representation of a set of frames used by four women's magazines in discussing breast cancer. The point of this presentation is to suggest which concepts, ideas, and media sources—here called “frames”—are similar, which cluster in small areas, and which are far apart. If all data points were equidistant from one another to begin with, there would be no point in scaling down the dimensionality of these data. Apparently, the standard against which MDS results become noteworthy is that of equal differences.

## Images, Portrayals, Semantic Nodes, and Profiles

## 10.5

Content analysts often focus on one or a few concepts, persons, or events and seek to ascertain how they are depicted or portrayed in text and what symbolic qualities readers might find associated with them. In the content analysis literature, titles like “Medicine in Medieval Literature,” “The Role of Scientists in Popular Media,” “The Human Body in Women's and Men's Magazines,” “How the Portrayal of the United States Has Shifted in Arab Dailies,” and “The Public Image of AT&T” abound. Researchers seek to answer questions of this kind by analyzing the linguistic or textual contexts in which references to the selected ideas occur.

In attribution analysis, the researcher tabulates the adjectives used to describe a chosen concept. A single list of attributes is quite uninformative, however, unless it is compared with some other list that provides a standard against which deviations may be noted. In a comparative attribution analysis, at least two lists are contrasted—for example, the image of one candidate for political office may be compared with the image of his or her opponent; or the portrayals of a country in textbooks before a war may be compared with those after that war; or the way one medium depicts a political scandal may be compared with how another medium covers that scandal. The analyst compares the lists to find out what attributes they share and what attributes distinguish among them. If all the attributes are shared among all the lists, there is little for the analyst to say. This reveals the standard that is common to this kind of analysis, the *sharing of attributions* against which differences in portrayals become noteworthy. Some researchers who conduct attribution analyses use expectations as a basis for comparison, reporting on how and how

much a given image deviates from the typical or usual. Unless a researcher has data on such expectations, formal tests may not be applicable. However, verbal highlights of what is unexpected or abnormal are common to many interpretations of images, portrayals, and the like (see the discussion of interactive-hermeneutic explorations in Chapter 11, section 11.6).

Another standard, common largely in linguistics, appears in the comparison of the linguistic context of one word or expression with the set of all grammatically and semantically acceptable contexts in which that word or expression can occur. The subset of actually used linguistic contexts is then equated with the meaning of the word or expression. This idea can easily be expanded to the meanings of politicians, professionals, academic disciplines, and countries.

Thus the notion of “attribute” should not be construed too narrowly. The image of a U.S. president that spin doctors and advertisers are so worried about can hardly be reduced to a list of adjectives. This would be a convenient but limited operationalization. It may have to include a president’s speeches, editorials discussing what the president does, opinion polls, even cartoons presenting that president’s public or private life. What is particular about the image of U.S. presidents is how what is said about them differs from what is said about comparable other personalities. Similarly, the image of, say, human genetics in science fiction makes sense only in comparison with how other scientific theories enter this genre.

Computer-aided text analysis (CATA), which I discuss in depth in Chapter 11, has provided us with several useful devices for the analysis of images and portrayals. One is the KWIC (keyword in context) list, a tabulation of sentences or text fractions that contain a particular word or phrase. Figure 11.1 shows such a tabulation for the word *play*. Weber (1984, p. 131) compared the KWIC lists for the word *rights* as used by Republicans and Democrats and found significant differences in how the two groups employed the word; these differences are what make Weber’s findings interesting. (See Chapter 11 for a fuller discussion of Weber’s study.) Researchers can examine the contexts of keywords or key phrases by using the “find” function of ordinary word processing programs, although this is a cumbersome method. Qualitative text analysis software moves from listing the contexts of single words to listing the contexts of categories of textual units (see Chapter 11, section 11.6).

Analyzing the nodes of semantic networks in terms of how one node is connected to others follows the same logic. For example, Figure 11.5 depicts the concept “hacker” as it appears in the narratives of students describing each other. In such networks, nodes are typically characterized in one of two ways:

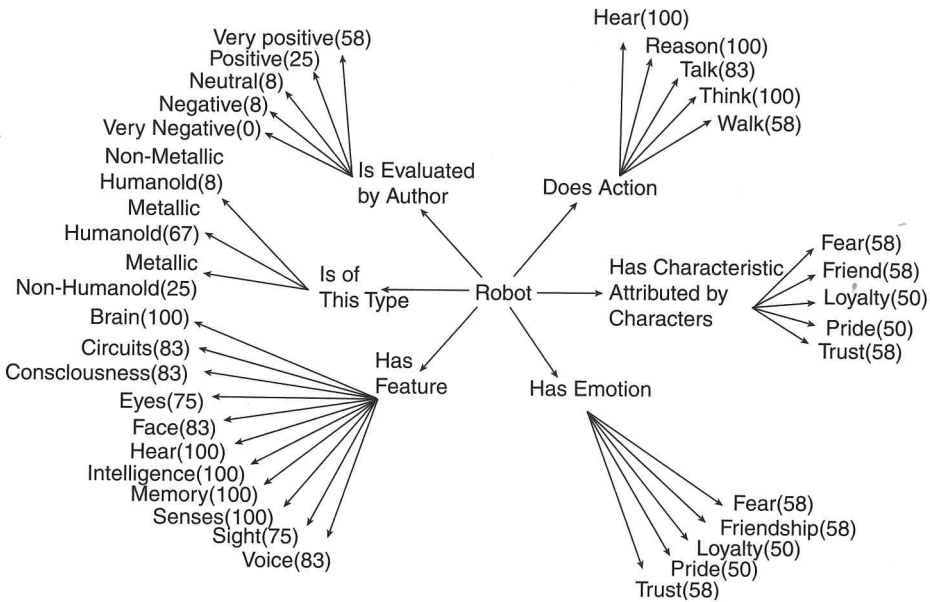
- They may be characterized in terms of measures that describe their position within a network—for example, with how many other nodes they are connected, their centrality or peripherality, or how often they occur. Carley (1997) has measured the positional properties of nodes in terms of density, conductivity, and intensity.
- They may be characterized in terms of the semantic connections between them and any other nodes. Figure 10.7, for example, depicts the semantic connections among nodes found by researchers who examined the characteristics attributed to robots in post-1960s texts (Palmquist, Carley, & Dale, 1997). This figure also displays the percentages of texts in which the connections occur.

Comparison of the linguistic environments in which a concept occurs gives rise to a variety of analytical possibilities. Two concepts that occur (or can occur) in the same linguistic environment are interchangeable, have the same meanings, and are considered synonymous. Concepts that mediate between many other concepts are the central concepts of a belief system, a story, or a discourse, which a network represents. An analysis of the environments that two concepts do not share elucidates differences in their meanings. “True” opposites share the environments of their genus but differ in everything else.

Figure 10.7 is one of several “maps” that Palmquist et al. (1997) compared in their examination of texts involving robots written before, during, and after 1960. What they found supported their hypotheses about changes in emotions associated with the robot image—over time, the texts showed emerging trust, loyalty, and friendship that increasingly counterbalanced persistent fears.

When analysts use profiles—whether of potential authors of unsigned documents, applicants for a certain job, or persons with mental illnesses—they apply the same interpretive standard, but with the addition that the attributes, correlates, or linguistic environments must be predictive of something. That is, they must answer a question, such as, Who wrote the unsigned document? Who is most likely to succeed in a given job? How should a therapist treat a patient who has a particular manner of talking?

Take the analysis of plagiarism as an example. Suppose that there are two literary works by different authors, A and B, and B is alleged to have plagiarized



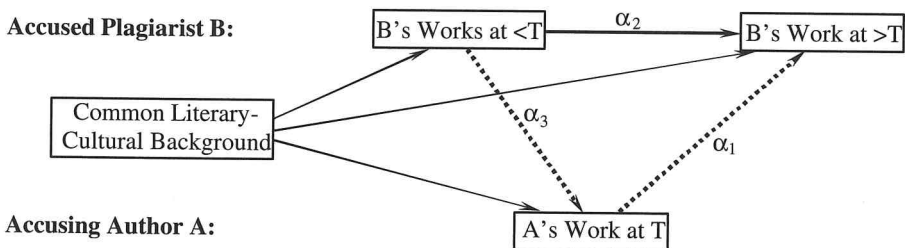
**Figure 10.7** The Robot Image: Robot Types, Features, and Actions From Post-1960s Texts

Source: Palmquist et al. (1997, p. 178, fig. 10.4).



A's work. Suspicions that one author has plagiarized the work of another are usually grounded in the recognition of surprising similarities between the two works in word choices, grammatical constructions, plot structure, outline, and so on. Even if the similarities are compelling, they do not constitute sufficient evidence of plagiarism. Before content analysts can enter a plagiarism dispute, it must be established that B had access to A's work before or while writing the disputed work. This is the easy part, to be addressed by a court of law. Once the accessibility of A's work to B has been established, analysts can focus their attention on how the similarities between the two works can be explained. Figure 10.8 diagrams the relationships that content analysts may have to consider. The similarity or agreement  $\alpha_1$  could be due to B's shameless copying of A's work, B's creativity (chance), or A's and B's common literary and/or cultural background. If the similarity  $\alpha_1$  can be shown to exceed  $\alpha_2$  substantially, this would add weight in favor of plagiarism on B's part. If  $\alpha_3$  exceeds  $\alpha_1$  to a degree better than chance, then A may actually be the plagiarist of B's previous work, rather than B having plagiarized A.

Authors, by definition, create new literature, and A and B could have come to these similarities on separate paths, especially if they are acquainted with each other's previously published work. Previous works may not be available for comparison or may not be considered relevant when the similarities being examined concern content, subject matter, or unusual personal experiences. But even the most imaginative writers rely on a background of literature, education, cultural practices, media exposure, and common sense that they share widely with others—otherwise their works would be unintelligible. This common background provides authors with a vocabulary of metaphors, sayings, myths, and themes that they weave into their writing. Most similarities between different authors' works are due to the background they share without realizing it. In a famous plagiarism case concerning a book about teaching in New York, the similarity turned out to be explained by the fact that, unknown to plaintiff and defendant, they had both taught in the same classroom in different years. If one subtracts the vocabulary and background that A and B share from the



**Figure 10.8** Comparisons of Works Needed to Establish or Dispel Accusations of Plagiarism

profiles of the two works, one is left with two profiles whose similarity or difference can be explained by creativity (or chance) or plagiarism. If the remaining similarities are well above chance, this finding might support a charge of plagiarism. The analysis of images, portrayals, semantic nodes, or profiles can lead in numerous directions.

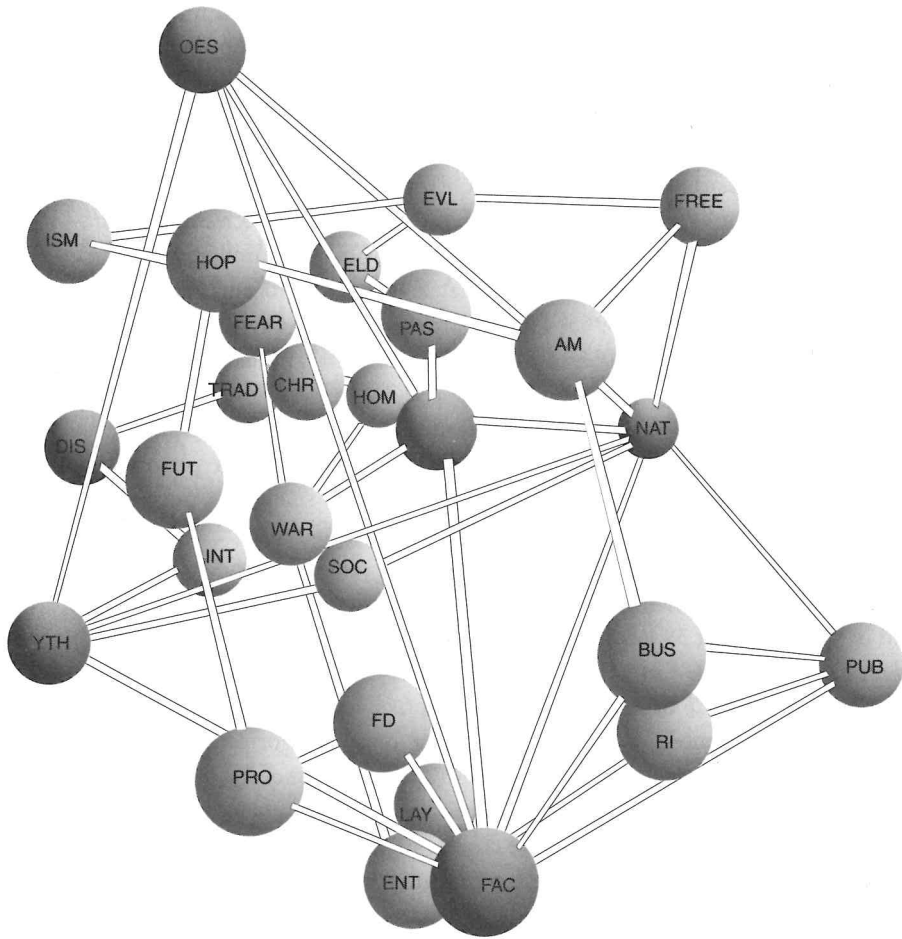
## Contingencies and Contingency Analysis 10.6

Contingency analysis is a technique that enables researchers to infer networks of associations from patterns of co-occurrences in texts, whether they are generated by a source or attended to by readers. Contingency analysis started with the observation that symbols often occur in pairs of opposites, that concepts or ideas form clusters. Contingency analysis is based on the assumption (analytical construct) that concepts that are closely associated cognitively will also be closely related proximally. Content analysts have successfully applied this assumption to individual authors, to social groups with common prejudices or ideological commitments, and to whole cultures permeated by cultural stereotypes or conventions. Experiments have shown not only that statistical contingencies in messages are a reflection of the associations in the mind of their authors, but that exposure to them can cause corresponding associations in their receivers as well, followed by the reproduction of these contingencies in speech, so that contingency analysis can be used to infer associations not only in the sources of texts but also in the audiences that are exposed to such statistical contingencies (Osgood, 1959, pp. 55–61). Regardless of these correlational validations (see Chapter 13, section 13.2.5), contingency analysis is an analytical technique in its own right.

Contingency analysis starts with a set of recording units, each of which is characterized by a set of attributes, concepts, or features that are either present or absent. The choice of recording units is important insofar as such units must contain sufficient numbers of co-occurrences. A word is too small a unit. A sentence usually contains several concepts, but units larger than sentences tend to be more productive. Osgood (1959), who first outlined this analysis, illustrated the steps involved in his analysis of 38 talks given by W. J. Cameron on the *Ford Sunday Evening Hour* radio program. First, Osgood regarded each talk as one recording unit and recorded the presence or absence of 27 conceptual categories in each talk. In the second step, he counted the co-occurrences of these categories and entered them in a square matrix of all pairs of categories (attributes, concepts, or features). In the third step, he tested the statistical significance of these co-occurrences. Co-occurrences that are significantly above chance suggest the presence of associations, whereas co-occurrences that are significantly below chance suggest the presence of dissociations.

The interpretive standard built into this technique is that of co-occurrences by chance, of course, indicating neither association nor dissociation and therefore not supporting inferences about the cognition of authors or readers. Osgood

plausibly argues that both directions of deviation from chance are of psychological importance. The association pattern that Osgood inferred from this rather small data set is depicted in Figure 10.9. Here, mentions of factories, industry, machines, production, and so on (FAC) tended to be associated with mentions of progress (PRO); Ford and Ford automobiles (FD); free enterprise and initiative (ENT); laymen, farmers, shopkeepers, and the like (LAY); and business, selling, and the like (BUS). But when Cameron used these concepts, he *avoided* talking about (to dissociate them from) such categories as youth (YTH), intellectuals, lily-livered bookmen, and so on (INT), and disease (DIS), which form another cluster of associations, dissociated from the former cluster. The figure shows also associations among violence and destruction (DES); assorted “isms,” such as communism,



**Figure 10.9** Spatial Representation of an Association Structure

Source: Osgood (1959, p. 68, fig. 4).

fascism, and totalitarianism (ISM); fear and bewilderment (FEAR); and sundry evils (EVL) (Osgood, 1959, pp. 67–68). Even without having heard these speeches, one can get a sense of the mentality of the speaker and of the times in which these speeches were broadcast.

The fundamental assumption underlying the analysis of contingencies is that co-occurrences in texts indicate associations in someone's mind or in existing cultural practices. This assumption, along with the idea of neuronal networks, has motivated Woelfel (1993, 1997) to develop software that allows a researcher to tabulate all co-occurrences of words within a sliding window of a specified length (e.g., 100 characters) and then compute clusters of contingencies. (I discuss this software, CatPac, in more detail in Chapter 11, section 11.4.2.) Incidentally, this idea underlies computational procedures that have been given the fancy name of "data mining." Text searches can identify contingencies as well within other kinds of windows, sentences, paragraphs, and documents. Thus several computer-aided approaches to text attend to contingencies within linguistic contexts. What distinguishes the latter from what Osgood had proposed is the absence of human readers, coders, or transcribers of the categories that are subjected to contingency analysis. Woelfel's aim is to bypass human readers or coders altogether, but the results that such software produces are bound to be shallow compared with analyses in which intelligent human readers are involved.

When tables of possible co-occurrences become very large, analysts may find it difficult to conceptualize the results. Examining a matrix of something like  $200 \times 200$  associations between concepts, which is not unusual in content analysis, is a formidable task, and analysts trying to discover patterns in such a flood of numerical data are likely to overlook important relationships. Then clustering becomes important.

## Clustering 10.7

Clustering operationalizes something humans do most naturally: forming perceptual wholes from things that are connected, belong together, or have common meanings, while separating them from things whose relationships seem accidental or meaningless. Clustering is closely allied with the conception of content as a representation, inviting abstraction, producing a hierarchy of representations that, on any one level, preserve what matters and omit only insignificant details from the original data. Procedurally, clustering either works from the bottom up, by lumping together objects, attributes, concepts, or people according to what they share, or proceeds from the top down, by dividing sets of such entities into classes whose boundaries reflect the more important differences between them. The direction that clustering takes results from the analyst's choices of the similarity measure and the clustering criterion. Clustering techniques differ widely regarding these. Contingency is but one similarity measure; others are agreement, correlation, proximity, the number of shared attributes, and common meanings, either by semantic definition or by relations within a thesaurus.

The choice of a clustering criterion is decisive for the kind of clusters a particular analysis provides. Some clustering criteria create long and snakelike clusters, whereas others produce compact and circular clusters. Some are sensitive to how much diversity accumulates within a cluster, others are not, assuring only the largest dissimilarities between the forming clusters (Krippendorff, 1980a). Under ideal circumstances, a clustering criterion reflects the way clusters are formed in the reality of the data source and relies on semantic similarities rather than purely syntactical ones. Content analysts must bear in mind that different clustering procedures may yield vastly different results; thus, to avoid ending up relying on arbitrary findings, they must always justify their use of particular clustering techniques in relation to the contexts of their analyses.

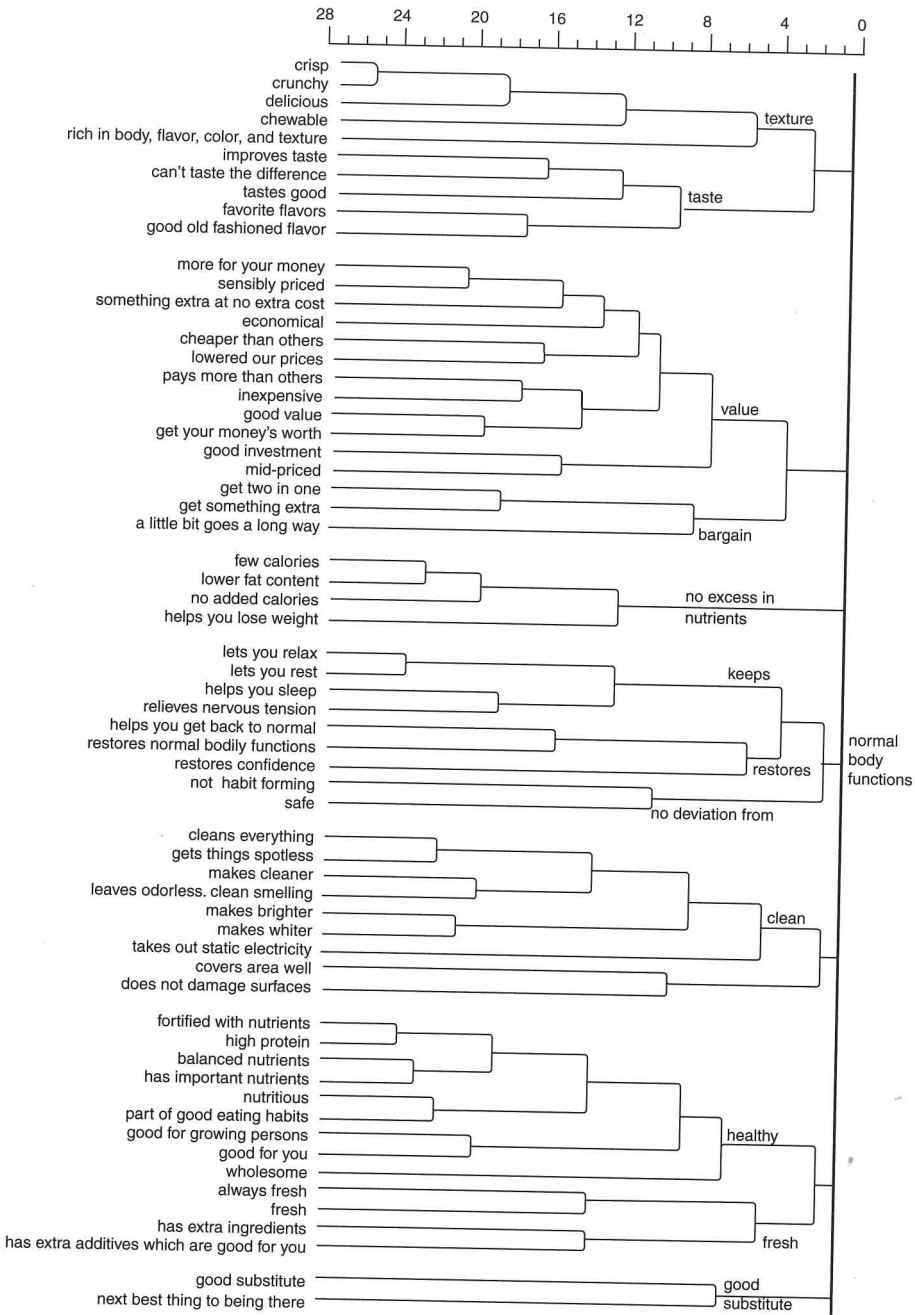
The most common of the available clustering procedures consists of the following iterative steps:

1. Within a matrix of similarity measures, search for two clusters (initially of two unclustered objects) that are, by the chosen criterion, most similar and the merger of which will least affect the overall measure of the differences in the data.
2. Lump these, taking into account the losses incurred within the newly formed cluster.
3. Recompute all measures of similarity with the newly formed cluster, thereby creating a new matrix of similarity measures within which the next two candidates for lumping are to be found.
4. Record the clustering step taken and the losses incurred for the user to retrace.
5. Repeat steps 1 through 4 until there is nothing left to merge (see Krippendorff, 1980a).

For a small amount of data and simple criteria, an analyst may even do clustering by hand.

Clustering steps are typically recorded in the form of so-called dendrograms, which are treelike diagrams that indicate when and which objects are merged and the losses the clustering incurred. Figure 10.10 shows a fraction of Dziurzynski's (1977) analysis of some 300 television advertising appeals. The resulting classification of appeals into those referring to texture, taste, value, and bargain appears to have considerable face validity.

As suggested above, clustering is popular in content analysis because, unlike factor analysis and multidimensional scaling, it is based on intuitively meaningful similarities among units of analysis, and its resulting hierarchies resemble the conceptualization of text on various levels of abstraction. This is why so many clustering algorithms are available. Often, however, the creators of these algorithms do not reveal how the algorithms work, and that puts the burden of proving their structural validity on the content analysts who use them (see Chapter 13, section 13.2.3).



**Figure 10.10** Part of a Large Dendrogram Resulting From Clustering Advertising Appeals

Source: Adapted from Dziurzynski (1977, pp. 25, 39, 40, 41, 50, figs. 3, 6, 11, 12, 13, 14, 53).