

BUILDING NON-NORMATIVE SYSTEMS - THE SEARCH FOR ROBUSTNESS
AN OVERVIEW

Mitchell P. Marcus
Bell Laboratories
Murray Hill, New Jersey 07974

Many natural language understanding systems behave much like the proverbial high school english teacher who simply fails to understand any utterance which doesn't conform to that teacher's inviolable standard of english usage. But while the teacher merely pretends not to understand, our systems really don't.

The teacher consciously stonewalls when confronted with non-standard usage to prescribe quite rigidly what is acceptable linguistic usage and what is not. What is so artificial about this behaviour, of course, is that our implicit linguistic models are descriptive and not prescriptive; they model what we expect, not what we demand. People are quite good at understanding language which they, when asked, would consider to be non-standard in some way or other.

Our programs, on the other hand, tend to be very rigid. They usually fail to degrade gracefully when their internal models of syntax, semantics or pragmatics are violated by user input. In essence, the models of linguistic well-formedness which these programs embody become normative; they prescribe quite rigidly what is considered standard linguistic usage and what isn't.

Old solutions to this problem include extending a system's linguistic coverage or intentionally excluding linguistic constraints that are occasionally violated by speakers. But neither of these approaches changes the fundamental situation - that when confronted with input which fails to conform to the system builder's expectations, however broad and however loose, the system will entirely reject the input. Furthermore, these techniques bar a system from utilizing the fact that people normally do obey certain linguistic standards, even if they violate them on occasion.

More recently, a range of approaches have been investigated that allow a system to behave more robustly when confronted with input which violates its designer's expectations about standard english usage. Most of this work has been within the realm of syntax. These techniques allow grammars to be descriptive without being normative. This panel focuses on these techniques for building what might be termed non-normative systems. Panelists were asked to consider the following range of issues:

Are there different kinds of non-standard usage? Candidates for a subclassification of non-standard usage might include the telegraphic language of messages and newspaper headlines; the informal colloquial use of language, even by speakers of the standard dialect; non-standard dialects; plain out-and-out grammatical errors; and the specialized sublanguage used by experts in a given domain. To what extent do these various forms have different properties, and are there independently characterizable dimensions along which they differ? What kinds of generalizations can be expressed about each of them individually or about non-standard usage in general?

What are the techniques for dealing with non-standard input robustly? A range of techniques have been discussed in the literature which can be invoked when a system is faced with input which is outside the subset of the language that its grammar describes. These include: (a) the use of special "un-grammatical" rules, which explicitly encode facts about non-standard usage; (b) the use of "meta-rules" to relax the constraints imposed by classes of rules of the grammar; (c) allowing flexible interaction between syntax and semantics, so that semantics can directly analyze substrings of syntactic fragments or individual words when full syntactic analysis fails. How well do these techniques, and others, work with respect to the dimensions of non-standard input discussed above? What are the relative strengths and weaknesses of each of these techniques?

To what extent are each of these techniques useful if one's goal is not to build a system which understands input, even if non-standard; but rather to build an explicitly normative system which can either (1) pinpoint grammatical errors, or (2) correct errors after pinpointing them? (Ironically, a system can be normative in a useful way only if it can understand what the user meant to say.)

Are there more general approaches to building systems that degrade gracefully that can be applied to this set of problems?

And finally, what the near- and long-term prospects for application of these techniques to practical working systems?