Natural Language Processing Semantic Aspects

Epaminondas Kapetanios Doina Tatar Christian Sacarea



Natural Language Processing

Semantic Aspects

Haldane: Where is the **bedeutung** of a proposition in your system, Turing? It is worth talking in terms of a universal grammar in mind, but it is also possible to construct meaningless propositions. **Turing:** For example?

Haldane: "Red thoughts walk peacefully".

..... This is an example of a proposition which is formulated correctly, according to the rules of the English grammar. If your theory is correct, then I was able to construct such a proposition, since I was able to activate the English version of the universal grammar in my mind, the semantic content of which equals to zero. Where can I find in your theory that this proposition makes no sense?

Turing: It is very simple, I do not know.

-from the book The Cambridge Quintett, John Casti, 1998

...however, the "system" is (as regards logic) a free play with symbols according to (logically) arbitrarily given rules of the game. All this applies as much (and in the same manner) to the thinking in daily life as to a more consciously and systematically constructed thinking in the sciences.

—Albert Einstein, On Remarks of B. Russell's Theory of Knowledge, Ideas and Opinions, New York, 1954

Natural Language Processing Semantic Aspects

Epaminondas Kapetanios

University of Westminster Faculty of Science and Technology London, UK

Doina Tatar

Babes-Bolyai University Faculty of Mathematics and Computer Science Cluj-Napoca, Romania

Christian Sacarea

Babes-Bolyai University Faculty of Mathematics and Computer Science Cluj-Napoca, Romania



CRC Press is an imprint of the Taylor & Francis Group, an **informa** business A SCIENCE PUBLISHERS BOOK CRC Press Taylor & Francis Group 6000 Broken Sound Parkway NW, Suite 300 Boca Raton, FL 33487-2742

© 2014 by Taylor & Francis Group, LLC CRC Press is an imprint of Taylor & Francis Group, an Informa business

No claim to original U.S. Government works Version Date: 20131021

International Standard Book Number-13: 978-1-4665-8497-6 (eBook - PDF)

This book contains information obtained from authentic and highly regarded sources. Reasonable efforts have been made to publish reliable data and information, but the author and publisher cannot assume responsibility for the validity of all materials or the consequences of their use. The authors and publishers have attempted to trace the copyright holders of all material reproduced in this publication and apologize to copyright holders if permission to publish in this form has not been obtained. If any copyright material has not been acknowledged please write and let us know so we may rectify in any future reprint.

Except as permitted under U.S. Copyright Law, no part of this book may be reprinted, reproduced, transmitted, or utilized in any form by any electronic, mechanical, or other means, now known or hereafter invented, including photocopying, microfilming, and recording, or in any information storage or retrieval system, without written permission from the publishers.

For permission to photocopy or use material electronically from this work, please access www.copyright.com (http://www.copyright.com/) or contact the Copyright Clearance Center, Inc. (CCC), 222 Rosewood Drive, Danvers, MA 01923, 978-750-8400. CCC is a not-for-profit organization that provides licenses and registration for a variety of users. For organizations that have been granted a photocopy license by the CCC, a separate system of payment has been arranged.

Trademark Notice: Product or corporate names may be trademarks or registered trademarks, and are used only for identification and explanation without intent to infringe.

Visit the Taylor & Francis Web site at http://www.taylorandfrancis.com

and the CRC Press Web site at http://www.crcpress.com

Preface

Communication has always played a pivotal role in the evolution of human culture, societies and civilisation. From symbols, cave paintings, petroglyphs, pictograms, ideograms and alphabet based forms of writing, to computerised forms of communication, such as the Web, search engines, email, mobile phones, VoIP Internet telephony, television and digital media, communication technologies have been evolving in tandem with shifts in political and economic systems.

Despite the emergence of a variety of communication means and technologies, natural language, or ordinary language, signed, written or spoken, remained the main means of communication among humans with the processing of natural language being an innate facility adhered to human intellect. With the rise of computers, however, natural language has been contrasted with artificial or constructed languages, such as Python, Java, C++, a computer programming language, or controlled languages for querying and search, in that natural languages contribute to the understanding of human intelligence. Nonetheless, the rise of social networks, e.g., Facebook, Twitter, did not replace natural language as one of the main means of information and communication in today's human civilisation.

What proved, however, to be an innate ability of humans from an early age to engage in speech repetition, language understanding and, therefore, so quickly acquire a spoken vocabulary from the pronunciation of words spoken around them, turned out to be a challenge for computing devices as symbol manipulators following a predefined set of instructions. Given the overwhelming amount of text based data and information generated and consumed daily, which is also boosted by the speedy pace of technological advances, effective solutions have been sought after, in accordance with many academic, industrial and scholar activities, in order to computationally improve natural language processing and understanding with the purpose to make meaningful information and communication stand out in an amalgamation of humans and machines.

Hence, it is not surprising that many books and research publications, around this topic, saw the light of this world amid all technological advances, which have been primarily geared towards faster communication rather than a qualitative one. The main tenor, what so ever, has been given through the lenses of text analytics, text mining, natural languages and information systems, information retrieval, as well as knowledge representation, artificial intelligence and machine learning. In all these contributions, a rather mathematical and statistical approach to natural language processing and understanding, than a pure linguistics based one, prevails, particularly when it comes to semantics and pragmatics based processing of text based data and information.

As having a long standing contribution to computational natural language processing and understanding, as well as to the underpinning mathematics, the authors epitomise, in this book, their experience and knowledge in a series of classic research areas at the cross-roads of natural language processing with information retrieval and search, text mining, knowledge representation, formal concept analysis and further mathematical aspects, with a particular emphasis on semantic aspects. To this extent, the book is, by no means, an exhaustive reference list of related work, or a handbook for natural language processing, however, it does provide a roadmap for those aspiring to contribute to world knowledge in the area through the lenses of semantic computing. Besides, the book aspires to guide all those academics, scholars and researchers, who wish to engage in this world knowledge contribution, through the major challenges and bumpy road ahead, as well as through a methodological baseline based on algorithmic and mathematical thinking, which underpins any serious attempt in computational approaches.

In this context, the first part of the book (Part I) introduces the reader into the main key challenges when it comes to representing and extracting meaning with such a symbol based system called natural language. The second part (Part II) discusses those mathematical aspects, which are considered fundamental for semantics based natural language processing. From a didactical point of view, some traditional mathematical concepts such as Lattice Theory, upon which Formal Concept Analysis is based, are being addressed in order to provide a second thought about the recurrent problems mainly caused by the shortly lived memory of classical studies in the field. Part III embarks on the knowledge representation aspects related with natural language processing in the flavour of measuring similarity among words, the pivotal role of semantics in query languages, as well as attempts to specify universal grammar for natural languages in the context of multilingual querying. Finally, part IV discusses knowledge extraction aspects related with natural language processing such as word sense disambiguation, text segmentation and summarisation, text entailment, named entity recognition.

> Epaminondas Kapetanios Doina Tatar Christian Sacarea

Contents

| Preface | v |
|---|--|
| PART I: Introduction | |
| The Nature of Language Syntax versus Semantics Meaning and Context The Symbol Grounding Problem | 3 3 8 13 |
| PART II: Mathematics | |
| 2. Relations 2.1 Operations with Relations 2.2 Homogenous Relations 2.3 Order Relations 2.4 Lattices. Complete Lattices 2.5 Graphical Representation of Ordered Sets 2.6 Closure Systems. Galois Connections | 19 22 25 32 34 36 38 |
| 3. Algebraic Structures 3.1 Functions 3.2 Binary Operations 3.3 Associative Operations. Semigroups 3.4 Neutral Elements. Monoids 3.5 Morphisms 3.6 Invertible Elements. Groups 3.7 Subgroups 3.8 Group Morphisms 3.9 Congruence Relations 3.10 Rings and Fields | $\begin{array}{c} \textbf{43} \\ 43 \\ 46 \\ 47 \\ 49 \\ 50 \\ 54 \\ 58 \\ 61 \\ 64 \\ 65 \end{array}$ |

| х | Natural | Language | Processing: | Semantic | Aspects |
|---|---------|----------|-------------|----------|---------|
|---|---------|----------|-------------|----------|---------|

| 4. | Linear Algebra 4.1 Vectors | 68 68 |
|----|--|-----------------|
| | 4.2 The space \mathbb{R}^n | 70 |
| | 4.3 Vector Spaces Over Arbitrary Fields | 72 |
| | 4.4 Linear and Affine Subspaces | $\overline{74}$ |
| | 4.5 Linearly Independent Vectors. Generator Systems. Basis | 5 79 |
| | 4.5.1 Every vector space has a basis | 83 |
| | 4.5.2 Algorithm for computing the basis of a generated sub-space | 90 |
| 5. | Conceptual Knowledge Processing and Formal | 94 |
| | Concept Analysis | |
| | 5.1 Introduction | 94 |
| | 5.2 Context and Concept | 96 |
| | 5.3 Many-valued Contexts | 106 |
| | 5.4 Finding all Concepts | 107 |
| | PART III: Knowledge Representation for NLP | |
| 6. | Measuring Word Meaning Similarity | 121 |
| | 6.1 Introduction | 121 |
| | 6.2 Baseline Methods and Algorithms | 122 |
| | 6.2.1 Intertwining space models and metrics | 122 |
| | 6.2.2 Measuring similarity | 126 |
| | 6.3 Summary and Main Conclusions | 128 |
| 7. | Semantics and Query Languages | 129 |
| | 7.1 Introduction | 129 |
| | 7.2 Baseline Methods and Algorithms | 131 |
| | 7.2.1 The methodology | 131 |
| | 7.2.2 The theory on semantics | 133 |
| | 7.2.3 Automata theory and (query) languages | 137 |
| | 7.2.4 Exemplary algorithms and data structures | 145 |
| | 7.3 Summary and Major Conclusions | 160 |
| 8. | Multi-Lingual Querying and Parametric Theory | 162 |
| | 8.1 Introduction | 162 |
| | 8.2 Baseline Methods and Algorithms | 164 |
| | 8.2.1 Background theory | 164 |
| | 8.2.2 An example | 168 |
| | 8.2.3 An indicative approach | 170 |
| | 8.2.4 An indicative system architecture and | 174 |
| | 8.3 Summary and Major Conclusions | 179 |
| | 0.0 Summary and Major Conclusions | 110 |

PART IV: Knowledge Extraction and Engineering for NLP

| 9. | Word Sense Disambiguation | 183 | | | | | | | |
|-----|--|-----|--|--|--|--|--|--|--|
| | 9.1 Introduction | | | | | | | | |
| | 9.1.1 Meaning and context | | | | | | | | |
| | 2 Methods and Algorithms: Vectorial Methods in WSD | | | | | | | | |
| | 9.2.1 Associating vectors to the contexts | | | | | | | | |
| | 9.2.2 Measures of similarity | 188 | | | | | | | |
| | 9.2.3 Supervised learning of WSD by vectorial methods | 189 | | | | | | | |
| | 9.2.4 Unsupervised approach. Clustering contexts by | 190 | | | | | | | |
| | vectorial method | | | | | | | | |
| | 9.3 Methods and Algorithms: Non-vectorial Methods in WSD | 192 | | | | | | | |
| | 9.3.1 Naive Bayes classifier approach to WSD | 192 | | | | | | | |
| | 9.4 Methods and Algorithms: Bootstrapping Approach | 193 | | | | | | | |
| | of WSD | | | | | | | | |
| | 9.5 Methods and Algorithms: Dictionary-based | 196 | | | | | | | |
| | Disambiguation | | | | | | | | |
| | 9.5.1 Lesk's algorithms | 196 | | | | | | | |
| | 9.5.2 Yarowsky's bootstrapping algorithm | 197 | | | | | | | |
| | 9.5.3 WordNet-based methods | 198 | | | | | | | |
| | 9.6 Evaluation of WSD Task | | | | | | | | |
| | 9.6.1 The benefits of WSD | | | | | | | | |
| | 9.7 Conclusions and Recent Research | | | | | | | | |
| 10. | Text Entailment | | | | | | | | |
| | 10.1 Introduction | 213 | | | | | | | |
| | 10.2 Methods and Algorithms: A Survey of RTE-1 and RTE-2 | 214 | | | | | | | |
| | 10.2.1 Logical aspect of TE | 216 | | | | | | | |
| | 10.2.2 Logical approaches in RTE-1 and RTE-2 | 218 | | | | | | | |
| | 10.2.3 The directional character of the entailment | 218 | | | | | | | |
| | relation and some directional methods in | | | | | | | | |
| | RTE-1 and RTE-2 . | | | | | | | | |
| | 10.2.4 Text entailment recognition by similarities | 220 | | | | | | | |
| | between words and texts | | | | | | | | |
| | 10.2.5 A few words about RTE-3 and the last RTE | 223 | | | | | | | |
| | challenges | | | | | | | | |
| | 10.3 Proposal for Direct Comparison Criterion | 223 | | | | | | | |
| | 10.3.1 Lexical refutation | 224 | | | | | | | |
| | 10.3.2 Directional similarity of texts and the | 227 | | | | | | | |
| | comparison criterion | | | | | | | | |

| | | 10.3.3 Two more examples of the comparison criterion | 228 |
|-----|--------|--|-----|
| | 10.4 | Conclusions and Recent Research | 229 |
| 11. | Text | Segmentation | 231 |
| | 11.1 | Introduction | 231 |
| | | 11.1.1 Topic segmentation | 232 |
| | 11.2 | Methods and Algorithms | 233 |
| | | 11.2.1 Discourse structure and hierarchical segmentation | 233 |
| | | 11.2.2 Linear segmentation | 236 |
| | | 11.2.3 Linear segmentation by Lexical Chains | 244 |
| | | 11.2.4 Linear segmentation by FCA | 248 |
| | 11.3 | Evaluation | 256 |
| | 11.4 | Conclusions and Recent Research | 260 |
| 12. | Text | Summarization | 262 |
| | 12.1 | Introduction | 262 |
| | 12.2 | Methods and Algorithms | 267 |
| | | 12.2.1 Summarization starting from linear | 267 |
| | | segmentation | |
| | | 12.2.2 Summarization by Lexical Chains (LCs) | 271 |
| | | 12.2.3 Methods based on discourse structure | 274 |
| | | 12.2.4 Summarization by FCA | 275 |
| | | 12.2.5 Summarization by sentence clustering | 280 |
| | 10.0 | 12.2.6 Other approaches | 283 |
| | 12.3 | Final Angle Summarization | 287 |
| | 12.4 | Evaluation 12.4.1. Conferences and Corners | 291 |
| | 195 | Conclusions and Recent Research | 294 |
| 10 | 12.0 | | 200 |
| 13. | Nam | led Entity Recognition | 297 |
| | 13.1 | Introduction | 297 |
| | 13.2 | Baseline Methods and Algorithms | 298 |
| | | 12.2.2. Machine learning techniques | 298 |
| | 13.3 | Summary and Main Conclusions | 303 |
| Bib | liogra | aphy | 311 |
| | | -E | 991 |
| ind | ex | | 33I |

PART I Introduction

CHAPTER 1

The Nature of Language

1.1 Syntax versus Semantics

It has been claimed many times in the past that humans are, somehow, born for grammar and speech as an innate ability to see the structure underlying a string of symbols. A classic example is the ease with which children pick up languages, which, in turn, has undergone evolutionary pressure. Without language, knowledge cannot be passed on, but only demonstrated. For instance, chimps can show the offsprings processes but cannot tell their about them, since a demonstration is required. Languages are, therefore, brought into connection with information, sometimes quite crucial. Language can help you to make plans. Many of the Spanish conquistadores who conquered Mesoamericans could not read, but their priests could. Moreover, being able to record language provides access to thousands of years of knowledge.

Generation and recognition of sentences pose two main problems for the concept of language as an assembly of a set of valid sentences. Though most textbooks deal with the understanding of the recognition of languages, one cannot ignore understanding the generation of language, if we aspire to understand recognition seriously.

A language can be described as a series of simple syntactic rules. For instance, English is a language defined with some simple rules, which are more loose than strict. This fact, however, may also highlight that a language can be hard to define with a only series of simple syntactic rules. Let us assume the following sentences:

- · John gesticulates
- · John gesticulates vigorously
- · The dog ate steak
- The dog ate ravenously

4 Natural Language Processing: Semantic Aspects

There are semantic rules (rules related to the meanings of sentences) in addition to the syntactic rules (rules regarding grammar). The rule usually specified, are strictly syntactic and, at least for computer languages, the easiest to formulate. The semantic rules, however, are notoriously difficult to formulate and are anchored in one's brain subconsciously, associating concepts with words and structuring the words into phrases and groups of phrases, which convey the meanings intended.

At a syntactic level and working towards some grammatical patterns or rules in English, one might be doing this consciously. There will always be a person or thing (a *subject*) and a verb describing an action (a *verb phrase*) in almost every language. In addition, there will sometimes be an *object* that the subject acts upon. In order to reflect on these abstract structures, one might find oneself using some other symbols acting as containers or patterns for sentences with a similar structure. For instance, one may end up with something like:

- *Subject* gesticulates
- *Subject* gesticulates vigorously
- *Subject* ate steak
- *Subject* ate ravenously

Next, abstract the verb phrases:

- Subject VerbPhrase
- Subject VerbPhrase
- Subject VerbPhrase steak
- Subject VerbPhrase

Finally, abstracting away the objects, we may end up with something like:

- Subject VerbPhrase
- Subject VerbPhrase
- Subject VerbPhrase Object
- Subject VerbPhrase

It is now easy to spot two main types of sentences that underpin the lexical-syntactic meanings of these four sentences:

- Subject VerbPhrase
- Subject VerbPhrase Object

You may also break down subject or verb phrases by having emerging sub-structures such as *noun* (*e.g.*, *John*) and *determinernoun* (*e.g.*, *The dog*) for subject phrases, or *verb* (e.g., ate) and *verbadverb* (e.g., ate ravenously) for verb phrases. Subsequently, you may end up with a finite language defined by the following rules of grammar:

- 1. A **sentence** is a subject followed by a verb phrase, optionally followed by an object.
- 2. A **subject** is a noun optionally preceded by a determiner.
- 3. A **verb** phrase is a verb optionally followed by an adverb.
- 4. A **noun** is John or dog.
- 5. A verb is gesticulates or ate.
- 6. An **adverb** is vigorously or ravenously.
- 7. An **object** is steak.
- 8. A **determiner** is The.

Despite the fact that the structure of these rules might seem to be right, here is exactly where the problem lies with the meaning of the sentences and the semantic rules associated with it. For example, the rules may allow you to say, "dog gesticulates ravenously," which is perfectly meaningless and a situation, which is frequently encountered as specifying grammars.

Having taken a look at how easily things might become quite complex when we need to define semantic rules on top of the syntactic ones, even with such a finite language, one can imagine that defining a strict grammar, i.e., including semantic rules, is almost impossible. For instance, a book on English grammar can easily become four inches thick. Besides, a natural language such as English is a moving target. For example, consider the difference between Elizabethan English and Modern English.

Again, as one discovers meta-languages in the next section one can bear in mind, that there is sometimes a gap between the language one means and the language one can easily specify. Another lesson learned is that using a language like English, which is neither precise nor terse, to describe other languages and, therefore, use it as a meta-language, one will end up with a meta-language with the same drawbacks.

The manner in which computer scientists have specified languages has been quite similar and is continuously evolving. Regardless of the variety and diversity of computer languages, semantic rules have rarely been an integral part of the language specification, if at all. They are mostly syntactic rules, which dominate the language specification. Take, for instance, context-free grammar (CFG), the first meta-language used extensively and preferred by most computer scientists. CFG specifications provide a list of rules with left and right hand sides separated by a right-arrow symbol. One of the rules is identified as the *start rule* or *start symbol*, implying that the overall structure of any sentence in the language is described by that rule. The left-hand side specifies the name of the substructure one is defining and the right hand side specifies the actual structure (sometimes called a *production*): a sequence of references to other rules and/or words in the language vocabulary.

Despite the fact that language theorists love CFG notation, most language reference guides use BNF (Backus- Naur Form) notation, which is really just a more readable version of CFG notation. In BNF, all rule names are surrounded by <...> and Æ is replaced with "::=". Also, alternative productions are separated by '|' rather than repeating the name of the rule on the left-hand side. BNF is more verbose, but has the advantage that one can write meaningful names of rules and is not constrained vis- à-vis capitalization. Rules in BNF take the form:

```
<rulename> ::= production 1
| production 2
...
| production n
```

Using BNF, one can write the eight rules used previously in this chapter as follows:

```
<Sentence> ::= <Subject> <VerbPhrase> <Object>
<Subject> ::= <Determiner> <Noun>
<VerbPhrase> ::= <Verb> <Adverb>
<Noun> ::= John | dog
<Verb> ::= gesticulates | ate
<Adverb> ::= vigorously | ravenously |
<Object> ::= steak |
<Determiner> ::= The |
```

Even if one uses alternatives such as YACC, the de facto standard for around 20 years, or ANTLR, or many other extended BNF (EBNF) forms, the highest level of semantic rule based specification one might achieve would be by introducing grammatical categories such as DETERMINER, NOUN, VERB, ADVERB and by having words such as *The, dog, ate, ravenously*, respectively, belonging to one of these categories. The intended grammar may take the following form,

sentence : subject verbPhrase (object)?; subject : (DETERMINER)? NOUN; verbPhrase : VERB (ADVERB)?; object : NOUN;

which still leaves plenty of space for construction of meaningless sentences such as *The dog gestured ravenously*. It is also worth mentioning that even with alternatives for CFG such as *regular expressions*, which were meant to simplify things by working with characters only and no rules referencing other ones on the righthand side, things did not improve towards embedding of semantic rules in a language specification. In fact, things turned out to be more complex with *regular expressions*, since without recursion (no stack), one cannot specify repeated structures.

In short, one needs to think about the difference between sequence of words in a sentence and what really dictates the validity of sentences. Even with the programming expression, if one is about to design state machinery capable of recognizing semantically sensitive sentences, the key idea must be that a sentence is not merely a cleverly combined sequence of words, but rather groups of words and groups of groups of words. Even with the programming expression (a[i+3)], humans can immediately recognize that there is something wrong with the expression, whereas it is notoriously difficult to design state machinery recognizing the faulty expression, since the number of left parentheses and brackets matches the number of one on the right.

In other words, sentences have a *structure* like this book. This book is organized into a series of chapters each containing sections, which, in turn, contain subsections and so on. Nested structures abound in computer science too. For example, in an object-oriented class library, classes group all elements beneath them in the hierarchy into categories (any kind of cat might be a subclass of feline etc...). The first hint of a solution to the underpowered state machinery is now apparent. Just as a class library is not an unstructured category of classes, a sentence is not just a flat list of words. Can one, therefore, argue that the role one gave each word, plays an equally large part in one's understanding of a sentence? Certainly, it does, but it is not enough. The examples used earlier, highlight the fact that structure imparts meaning very clearly. It is not purely the words though, nor the sequence that impart meaning. Can it also be argued that if state machines can generate invalid sentences, they must be having trouble with structure? These questions will be left unanswered for the time being, or perhaps in the near future. It turns out that even if we manage to define state machinery to cope with structure in a sentence, claiming that once semantic rules are perfectly defined, it is far reaching, since there are more significant issues to consider, as we will see in the following sections.

1.2 Meaning and Context

The difficulty of defining semantic rules to cope with meaningful states, operations or statements is exacerbated by the conclusions drawn from the study of "Meaning" as a key concept for understanding a variety of processes in living systems. It turns out that "Meaning" has an elusive nature and "subjective" appearance, which is, perhaps the reason why it has been ignored by information science. Attempts have been made to circumscribe a theory of meaning in order to determine the meaning of an indeterminate sign. Meaning-making, however, has been considered as the procedure for extracting the information conveyed by a message, in which the former is considered to be the set of values one might assign to an indeterminate signal. In this context, meaning-making is described in terms of a constraint-satisfaction problem that relies heavily on contextual cues and inferences.

The lack of any formalization of the concepts "meaning" and "context", for the working scientist, is probably due to the theoretical obscurity of concepts associated with the axis of semiotics in information processing and science. Even with regard to information and information flow, it has been argued that "the formulation of a precise, qualitative conception of information and a theory of the transmission of information has proved elusive, despite the many other successes of computer science" (Barwise and Seligman 1993). Since Barwise's publication, little has changed. Researchers in various fields still find it convenient to conceptualize the data in terms of information theory. By doing so, they are excluding the more problematic concept of meaning from the analysis. It is clear, however, that the meaning of a message cannot be reduced to the information content.

In a certain sense, the failure to reduce meaning to information content is like the failure to measure organization through information content. Moreover, the relevance of information theory is criticized by those who argue that when we study a living system, as opposed to artificial devices, our focus should be on meaning-making rather than information processing per se. In the context of artificial devices, the probabilistic sense of information prevails. Meaning, however, is a key concept for understanding a variety of processes in living systems, from recognition capacity of the immune system to the neurology of the perception. Take, for instance, the use of information theory in biology, as stated by (Emmeche and Hoffmeyer 1991). They argue that unpredictable events are an essential part of life and it is impossible to assign distinct probabilities to any event and conceptualize the behavior of living systems in terms of information theory. Therefore, biological information must embrace the "semantic openness" that is evident, for example, in human communication.

In a nutshell, it is worth mentioning that meaning has taken both main views, divorced from information and non-reducible to each other. The first is due to the fact that the concept of information relies heavily on "information theory" like Shannon's statistical definition of information, whereas the latter is due to the conception that information can broadly be considered as something conveyed by a message in order to provoke a response (Bateson 2000). Hence, the message can be considered as a portion of the world that comes to the attention of a cogitative system, human or non-human. Simply stated, information is a differentiated portion of reality (i.e., a message), a bit of information as a difference, which makes a difference, a piece of the world that comes to notice and results in some response (i.e., meaning). In this sense, information is interactive. It is something that exists in between the responding system and the differentiated environment, external or internal. An example, if one leaves one's house to take a walk, notices that the sky is getting cloudy, one is likely to change one's plans in order to avoid the rain. In this care the cloudy sky may be considered the message (i.e., the difference) and one's avoidance will be the information conveyed by the message (i.e., a difference that makes a difference). In this context, information and meaning are considered synonymous and without any clear difference between them. Though they are intimately related, they cannot be reduced to each other.

10 Natural Language Processing: Semantic Aspects

In the same spirit, Bateson presents the idea that a differentiated unit, e.g., *a word*, has meaning only on a higher level of logical organization, e.g., *the sentence*, only in context and as a result of interaction between the organism and the environment. In this sense, the internal structure of the message is of no use in understanding the meaning of the message.

The pattern(s) into which the sign is woven and the interaction in which it is located is what turns a differentiated portion of the world into a response by the organism. This idea implies that turning a signal (i.e., a difference) into a meaningful event (i.e., a difference that makes a difference) involves an active extraction of information from the message. Based on the suggestions, the following ideas have been suggested:

- a) Meaning-making is a procedure for extracting the information conveyed by a message.
- b) Information is the value one may assign to an indeterminate signal (i.e., a sign).

These ideas are very much in line with conceptions that see meaning-making as an active process that is a condition for information-processing rather than the product of informationprocessing per se. The most interesting things in the conception of meaning-making as an active process, are the three organizing concepts of a) indeterminacy of the signal, b) contextualization, c) transgradience. The indeterminacy (or variability) of the signal is an important aspect of any meaning-making process. *It answers the question what is the indeterminacy of the signal and why is it important for a theory of meaning-making*? The main idea is that in itself every sign/unit is devoid of meaning until it is contextualized in a higher-order form of organization such as a *sentence*. It can be assigned a range of values and interpretations.

For instance, in natural language the sign "shoot" can be used in one context to express an order to a soldier to fire his gun and in a different context as a synonym for "speak". In immunology, the meaning of a molecule's being an antigen is not encapsulated in the molecule itself. That is, at the most basic level of analysis, a sign has the potential to mean different things (i.e., to trigger different responses) in different contexts, a property that is known in linguistics as polysemy and endows language with enormous flexibility and cognitive economy. In the field of linguistics it is called pragmatics, which deals with meaning in context, the single most obvious way in which the relation between language and context is reflected in the structure of languages themselves is often called *deixis* (pointing or indicating in Greek). To this extent, linguistic variables (e.g., this, he, that) are used to indicate something in a particular context. They are indeterminate signals.

Nevertheless, the indeterminacy of a signal or word can be conceived as a constraint satisfaction problem. This, in turn, is defined as a triple {V, D, C}, where: (a) V is a set of variables, (b) D is a domain of values, and (c) C is a set of constraints {C1,C2, . . .,Cq}. In the context of semiotics, V is considered to be the set of indeterminate signals and D the finite set of interpretations/values one assigns to them. Based on the above definition, a sign is indeterminate if assigning it a value is a constraint-satisfaction problem. One should note that solving the constraint-satisfaction problem is a meaningmaking process, since it involves the extraction of the information conveyed by a message (e.g., to whom does the "he" refer?). However, rather than a simple mapping from V to D, this process also involves contextualization and inference.

The problematic notion of *context*, in the conception of meaningmaking as an active process, can be introduced better as an *environmental setting composed of communicating units and their relation in time and space*. The general idea and situation theory (Seligman and Moss 1997) is one possible way of looking into these aspects. In situation theory, a situation is *"individuals in relations having properties and standing in relations to various spatiotemporal relations"*. In a more general way, we can define a situation as a pattern, a meaning, an ordered array of objects that have an identified and repeatable form. In an abstract sense, a contextualization process can be conceived as a functor or a structure-preserving mapping of the particularities or the token of the specific occurrence onto the generalities of the pattern.

Regarding the interpretation of things as a constraints satisfaction problem, a context forms the constraints for the possible values (i.e., interpretations) that one may attribute to a sign. According to this logic, a situation type is a structure or pattern of situations. In other words, a situation is defined as a set of objects organized both spatially and temporally in a given relation. If this relation characterizes a set of situations, one can consider it a structure or a situation type. For example, the structure of hierarchical relations is the same no matter who the boss is. The situation type is one of hierarchical relations. Based on this type of a situation, we can make inferences about other situations. For example, violations of a rigid hierarchical relationship by a subordinate are usually responded to with another situation of penalties imposed by the superiors.

Although a sign, like the meaning of a word in a sentence, is indeterminate, in a given context one would like to use it to communicate only one meaning (i.e., to invite only one specific response) and not others. Therefore, the word disambiguation problem arises. In a sense, inferences via contextualization work as a zoom-in, zoom-out function. Contextualization offers the ability to zoom out and reflexively zoom back, in a way that constrains the possible values we may assign to the indeterminate signal. In other words, in order to determine the meaning of a microelement and extract the information it conveys, one has to situate it on a level of higher order of organization.

Let us consider the following example: "I hate turkeys". The vertical zooming-out from the point representing "I" to the point representing "human" captures the most basic denotation of "I," given that *denotation* is the initial meaning captured by a sign. As such, it could be considered the most reasonable starting point for contextualization. It is also a reasonable starting point because both evolutionary and ontological denotation is the most basic semiotic category. According to the *Oxford English Dictionary*, a turkey can be zoomed out to its closest ontological category, a "bird". This ontological category commonly describes any feathered vertebrate animal. Therefore, if we are looking for a function from the indeterminate sign "love" to a possible value/interpretation, the first constraint is that it is a relation between a *human being* and an *animal*.

There is, however, one more contextual cue. This is dictated by the denotation of dinner as a token of a *meal*. In that situation, there is a relationship of eating between *human beings* and *food*. Given that the zoomed-out concept of human beings for the sign "T" does participate in this relationship as well, another candidate value for the interpretation of "love" will arise, which, apparently, makes more sense, since it may be much closer to the meaning of the sentence "I hate turkeys". The situation where humans consume turkeys as food is the one giving meaning to this sentence.

Contextualization, however, is not sufficient to meaning-making processes. *Transgradience*, as the third dimension of the meaningmaking process, refers to the need for interpretation, inference and integration as a process in which *inferences are applied to a signal-in-context in order to achieve a global, integrated view of a situation*. In general terms, transgradience refers to the ability of the system to achieve a global view of a situation by a variety of means. An interesting parallel can be found by the immune system deciding whether a molecule is an antigen by means of a complex network of immune agents that communicate, make inferences, and integrate the information they receive. Further sub-dimensions may arise though, which could potentially complicate things: (1) the spatiotemporal array in which the situation takes place; (2) our background knowledge of the situation and (3) our beliefs concerning the situation.

In brief, our ability to extract the information that a word may convey as a signal, is a meaning-making process that relies heavily on contextual cues and inferences. Now the challenge is to pick the right situation, which constrains the interpretation values to be allocated to indeterminate signals, i.e., ambiguous words. Although one's understanding of semiotic systems has advanced (Sebeok and Danesi 2000) and computational tools have reached a high level of sophistication, one still does not have a satisfactory answer to the question of how the meaning emerges in a particular context.

1.3 The Symbol Grounding Problem

The whole discussion, so far, is underpinned by the assumption that the adherence of meaning to symbols and signals is a result of a meaning-making process rather than something intrinsic to the symbols and the chosen symbolic system itself. It is this innate feature of any symbolic system, which poses limitations to any symbol manipulator, to the extent to which one can interpret symbols as having meaning systematically. This turns interpretation of any symbol such as letters or words in a book parasitic, since they derive their meaning from us similarly, none of the symbolic systems can be used as a cognitive model and therefore, cognition cannot just be a manipulation of a symbol. Spreading this limitation would mean grounding every symbol in a symbolic system with its meaning and not leaving interpretation merely to its shape. This has been referred to in the 90s as the famous 'symbol grounding problem' (Harnad 1990) by raising the following questions : • "How can the semantic interpretation of a formal symbol system be made intrinsic to the system, rather than remain parasitic, depending solely on the meanings in our heads?"

• "How can the meanings of the meaningless symbol tokens, manipulated solely on the basis of their (arbitrary) shapes, be grounded in anything but other meaningless symbols?"

The problem of constructing a symbol manipulator able to understand the extrinsic meaning of symbols, has been brought into analogy with another famous problem of trying to learn Chinese from a Chinese/Chinese dictionary. It also sparked off the discussion about symbolists, symbolic Artificial Intelligence and the symbolic theory of mind, which has been challenged by Searle's "Chinese Room Argument". According to these trends, it has been assumed that if a system of symbols can generate indistinguishable behavior in a person, this system must have a mind. More specifically, according to the symbolic theory of mind, if a computer could pass the Turing Test in Chinese, i.e., if it could respond to all Chinese symbol strings it receives as input from Chinese symbol strings that are indistinguishable from the replies a real Chinese speaker would make (even if we keep on testing infinitely), the computer would understand the meaning of Chinese symbols in the same sense that one understands the meaning of English symbols. Like Searle's demonstration, this turns out to be impossible, for both humans and computers, since the meaning of the Chinese symbols is not intrinsic and depends on the shape of the chosen symbols. In other words, imagine that you try to learn Chinese with a Chinese/Chinese dictionary only. The trip through the dictionary would amount to a merry-go-round, passing endlessly from one meaningless symbol or symbol-string (the definientes) to another (the definienda), never stopping to explicate what anything meant.

The standard reply and approach of the symbolists and the symbolic theory of mind, which prevails in the views of semantic aspects in natural language processing within this book as well, is that the meaning of the symbols comes from connecting the symbol system to the world "in the right way." Though this view trivializes the symbol grounding problem and the meaning making process in a symbolic system, it also highlights the fact that if each definiens in a Chinese/Chinese dictionary were somehow connected to the world in the right way, we would hardly need the definienda. Therefore, this would alleviate the difficulty of picking out the objects, events and states of affairs in the world that symbols refer to.

With respect to these views, hybrid non-symbolic / symbolic systems have been proposed in which the elementary symbols are grounded in some kind of non-symbolic representations that pick out, from their proximal sensory projections, the distal object categories to which the elementary symbols refer. These groundings are driven by the insights of how humans can (1) discriminate, (2) manipulate, (3) identify and (4) describe the objects, events and states of affairs in the world they live in and they can also (5) "produce descriptions" and (6) "respond to descriptions" of those objects, events and states of affairs.

The attempted groundings are also based on discrimination and identification, as two complementary human skills. To be able to discriminate one has to judge whether two inputs are the same or different, and, if different, to what degree. Discrimination is a relative judgment, based on our capacity to tell things apart and discern the degree of similarity.

Identification is based on our capacity to tell whether a certain input is a member of a category or not. Identification is also connected with the capacity to assign a unique response, e.g., a name, to a class of inputs. Therefore, the attempted groundings must rely on the answers to the question asking what kind or internal representation would be needed in order to be able to discriminate and identify. In this context, iconic representations have been proposed (Harnad 1990). For instance, in order to be able to discriminate and identify horses, we need horse icons. Discrimination is also made independent of identification in that one might be able to discriminate things without knowing what they are. According to the same theorists, icons alone are not sufficient to identify and categorize things in an underdetermined world full with infinity of potentially confusable categories. In order to identify, one must selectively reduce those to "invariant features" of the sensory projection that will reliably distinguish a member of a category from any non-members with which it could be confused. Hence, the output is named "categorical representation". In some cases these representations may be innate, but since evolution could hardly anticipate all the categories one may ever need or choose to identify, most of these features must be learned from experience. In a sense, the categorical representation of a horse is probably a learned one.

16 Natural Language Processing: Semantic Aspects

It must be noted, however, that both representations are still sensory and non- symbolic. The former are analogous copies of the sensory projection, preserving its "shape" faithfully. The latter are supposed to be icons that have been filtered selectively to preserve only some of the features of the shape of the sensory projection, which distinguish members of a category from non-members reliably. This sort of non-symbolic representation seems to differ from the symbolic theory of mind and currently known symbol manipulators such as conventional computers trying to cope with natural language processing and their semantic aspects.

Despite the interesting views emerging from the solution approaches to the symbol grounding problem, the symbol grounding scheme, as introduced above has one prominent gap: no mechanism has been suggested to explain how the all-important categorical representations can be formed. How does one find the invariant features of the sensory projection that make it possible to categorize and identify objects correctly? To this extent, connectionism, with its general pattern learning capability, seems to be one natural candidate to complement identification. In effect, the "connection" between the names and objects that give rise to their sensory projections and icons would be provided by connectionist networks. Icons, paired with feedback indicating their names, could be processed by a connectionist network that learns to identify icons correctly from the sample of confusable alternatives it has encountered. This can be done by adjusting the weights of the features and feature combinations that are reliably associated with the names in a way that may resolve the confusion. Nevertheless, the choice of names to categorize things is not free from extrinsic interpretation of things, since some symbols are still selected to describe categories.

PART II Mathematics

CHAPTER 2

Relations

The concept of a relation is fundamental in order to understand a broad range of mathematical phenomena. In natural language, *relation* is understood as *correspondence*, *connection*. We say that two objects are related if there is a common property linking them.

Definition 2.0.1. Consider the sets *A* and *B*. We call (*binary*) *relation* between the elements of *A* and *B* any subset $R \subseteq A \times B$. An element $a \in A$ is in relation *R* with an element $b \in B$ if and only if $(a, b) \in R$. An element $(a, b) \in R$ will be denoted by aRb.

Definition 2.0.2. If A_1, \ldots, A_n , $n \ge 2$ are sets, we call an *n*-ary relation any subset $R \subseteq A_1 \times \ldots \times A_n$. If n = 2, the relation R is called *binary*, if n = 3 it is called *ternary*. If $A_1 = A_2 = \ldots = A_n = A$, the relation R is called *homogenous*.

In the following, the presentation will be restricted only to binary relations.

Remark 1 The direct product $A \times B$ is defined as the set of all ordered pairs of elements of A and B, respectively:

$$A \times B := \{(a, b) \mid a \in A, b \in B\}.$$

Remark 2 If A and B are finite sets, we can represent relations as cross tables. The rows correspond to the elements of A, while the columns correspond to the elements of B. We represent the elements of R, i.e., $(a, b) \in R$, by a cross (X) in this table. Hence, the relation R is represented by a series of entries (crosses) in this table. If at the intersection of row a with column b there is no entry, it means that a and b are not related by R.

Example 2.0.1. Relations represented as cross-tables

(1) Let $A = \{a\}$. There are only two relations on A, the empty relation, $R = \emptyset$, and the total relation, $R = A \times A$.



(2) $A := \{1, 2\}, B := \{a, b\}$. Then, all relations $R \subseteq A \times B$ are described by

| | a | b | | a | b | | a | b | | a | b |
|---|---|---|---|---|---|---|-----|---|---|---|---|
| 1 | x | | 1 | | x | 1 | | | 1 | | |
| 2 | | | 2 | | | 2 | x | | 2 | Х | |
| | a | b | | a | b | | a | b | | a | b |
| 1 | X | X | 1 | | | 1 | X | | 1 | | X |
| 2 | | | 2 | X | X | 2 | X | | 2 | | X |
| | | | | | | | | | | | |
| | a | b | | a | b | | a | b | | a | b |
| 1 | Х | | 1 | | X | 1 | X | X | 1 | X | |
| 2 | | Х | 2 | X | | 2 | X | | 2 | X | X |
| | a | b | | a | b | | a | b | | a | b |
| 1 | | Х | 1 | X | X | 1 | X | X | 1 | | |
| 2 | Х | Х | 2 | | Х | 2 | X | X | 2 | | |

Example 2.0.2. In applications, A, B, and $R \subseteq A \times B$ are no longer abstract sets, they have a precise semantics, while the relation R represents certain correspondences between the elements of A and B. The following example describes some arithmetic properties of the first ten natural numbers:

| | even | odd | div.by 3 | div.by 5 | div.by 7 | prime | $x^2 + y^2$ | $x^2 - y^2$ |
|----|------|-----|----------|----------|----------|-------|-------------|-------------|
| 1 | | Х | | | | | | |
| 2 | Х | | | | | Х | Х | |
| 3 | | Х | Х | | | Х | | Х |
| 4 | Х | | | | | | | |
| 5 | | Х | | Х | | Х | | Х |
| 6 | х | | Х | | | | | |
| 7 | | Х | | | Х | Х | | |
| 8 | Х | | | | | | Х | Х |
| 9 | | Х | Х | | | | | |
| 10 | Х | | | Х | | | | |

Example 2.0.3. Other relations

- (1) The divisibility relation in \mathbb{Z} : $R := \{(m, n) \in \mathbb{Z}^2 \mid \exists k \in \mathbb{Z} : n = km\}$.
- (2) $R := \{(x, y) \in \mathbb{R}^2 \mid x^2 + y^2 = 1\}$. This relation consists of all points located on the circle centered in the origin and radius 1.
- (3) The equality relation on set A:

$$\Delta_A := \{ (x, x) \mid x \in A \}.$$

- (4) The equality relation in \mathbb{R} consists of all points located on the first bisecting line.
- (5) The universal relation on a set *A* expresses the fact that all elements of that set are related to each other:

$$\nabla_A := \{(x, y) \mid x, y \in A\}.$$

(6) The empty relation means that none of the elements of A and B are related:

$$R = \emptyset \subseteq A \times B.$$

(7) Let $A = B = \mathbb{Z}$ and *R* the divisibility relation on \mathbb{Z} :

$$R := \{ (x, y) \in \mathbb{Z} \times \mathbb{Z} \mid \exists k \in \mathbb{Z}. \ y = kx \}.$$

22 Natural Language Processing: Semantic Aspects

2.1 Operations with Relations

Let (A_1, B_1, R_1) and (A_2, B_2, R_2) be two binary relations. They are equal if and only if $A_1 = A_2$, $B_1 = B_2$, $R_1 = R_2$.

Definition 2.1.1. Let *A* and *B* be two sets, *R* and *S* relations on $A \times B$. Then *R* is *included* in *S* if $R \subseteq S$.

Definition 2.1.2. Let *R*, $S \subseteq A \times B$ be two relations on $A \times B$. The *intersection* of *R* and *S* is defined as the relation $R \cap S$ on $A \times B$.

Definition 2.1.3. Let $R, S \subseteq A \times B$ be two relations on $A \times B$. The *union* of R and S is defined as the relation $R \cup S$ on $A \times B$.

Definition 2.1.4. Let $R \subseteq A \times B$ be a relation on $A \times B$. The *complement* of *R* is defined as the relation CR on $A \times B$, where

$$\mathbb{C}R := \{ (a, b) \in A \times B \mid (a, b) \notin A \times B \}.$$

Remark 3 If R and S are relations on $A \times B$ then

(1) $a(R \cap S)b \Leftrightarrow aRb$ and aSb.

(2) $a(R \cup S)b \Leftrightarrow aRb$ or aSb.

(3) $a(\mathbb{C}R)b \Leftrightarrow (a, b) \in A \times B$ and $(a, b) \notin R$.

Definition 2.1.5. Let $R \subseteq A \times B$ and $S \subseteq C \times D$ be two relations. The *product* or *composition* of *R* and *S* is defined as the relation $S \circ R \subseteq A \times D$ by

 $S \circ R := \{(a, d) \in A \times D \mid \exists b \in B \cap C. (a, b) \in R \text{ and } (b, d) \in S\}.$

If $B \cap C = \emptyset$, then $S \circ R = \emptyset$.

Definition 2.1.6. *Let* $R \subseteq A \times B$ be a relation. The *inverse* of R is a relation $R^{-1} \subseteq B \times A$ defined by

$$R^{-1} := \{ (b, a) \in B \times A \mid (a, b) \in R \}.$$

Theorem 2.1.1. *Let* $R \subseteq A \times B$, $S \subseteq C \times D$, and $T \subseteq E \times F$ be relations.

Then the composition of relations is associative:

 $(T \circ S) \circ R = T \circ (S \circ R).$

Proof. Let $(a, f) \in (T \circ S) \circ R$. By the definition of the relational product, there exists $b \in B \cap C$ with $(a, b) \in R$ and $(b, f) \in T \circ S$. By the same definition, there exists $d \in D \cap E$ with $(b, d) \in S$ and $(d, f) \in T$. Now, $(a, b) \in R$ and $(b, d) \in S$ imply $(a, d) \in S \circ R$. Together with $(d, f) \in T$, this implies that $(a, f) \in T \circ (S \circ R)$. Hence,

$$(T \circ S) \circ R \subseteq T \circ (S \circ R).$$