

1 Information Extraction and Information Technology

With Rik De Busser

1.1 Defining Information Extraction

A company wants to track the general sentiments about its newly released product in Web blogs. Another company wants to use the news feeds it bought from a press agency to construct a detailed overview of all technological trends in the development of semiconductor technologies. The company also wants a timeline of all business transactions involved in this development. A space agency allows astronauts to query large amounts of technical documentation by means of natural language speech. A government is gathering data on a natural disaster and wants to urgently inform emergency services with a summary of the latest data available. An intelligence agency is investigating general trends in terrorist activities all over the world. They have a database of millions of news feeds, minutes and e-mails and want to use these to get a detailed overview of all terrorist events in a particular geographical region in the last five years. A legal scholar is interested in studying the decisions of judges in divorce settlements and the underlying criteria. He or she has thousands of court decisions at his disposal. A biomedical research group is investigating a new treatment and wants to know all possible ways in which a specific group of proteins can interact with other proteins and what the exact results of these interactions are. There are tens of thousands of articles, conference papers and technical reports to study.

The above examples have a number of elements in common: (1) The requests for information; (2) The answer to this request is usually present in unstructured data sources such as text and images; (3) But, it is impossible for humans to process all data because there is simply too much of it; And

(4) computers are not able to directly query for the target information because it is not stored in a structured format such as a database but in unstructured sources. *Information extraction (IE)* is the subdiscipline of artificial intelligence that tries to solve this kind of problems.

Traditionally, information extraction is associated with template based extraction of event information from natural language text, which was a popular task of the *Message Understanding Conferences* in the late eighties and nineties (Sundheim, 1992). MUC information extraction tasks started from a predefined set of templates, each containing specific information slots that encode event types relevant to a very specific subject domain – for instance, terrorism in Latin America – and used relatively straightforward pattern matching techniques to fill out these templates with specific instances of these events from a corpus of texts. Patterns in the form of a grammar or rules (e.g., in the form of regular expressions) were mapped on the text in order to identify the information.

MUC was the first large scale effort to boost research into automatic information extraction and it would define the research field for the decades to come. Even at the time of writing, information extraction is often associated with template based pattern matching techniques. Unsurprisingly, the MUC legacy still resounds very strongly in Riloff and Lorenzen's definition of information extraction:

IE systems extract domain-specific information from natural language text. The domain and types of information to be extracted must be defined in advance. IE systems often focus on object identification, such as references to people, places, companies, and physical objects. [...] Domain-specific extraction patterns (or something similar) are used to identify relevant information.

(Riloff and Lorenzen, 1999, p. 169)

This definition represents a traditional view on what information extraction is and it more or less captures what this discipline is about: The extraction of information that is semantically defined from a text, using a set of extraction rules that are tailored to a very specific domain. The main points expressed by this definition are that an information extraction system identifies information in text, i.e., in an unstructured information source, and the information that adheres to predefined semantics (e.g., people, places etc.). However, we will see in the rest of the book that at present the scope of Riloff and Lorenzen's definition has become too limited. Information

extraction is not necessarily domain specific. In practice, the domain of the information to be extracted is often determined in advance, but this has more to do with technological limitations of the present state of the art than with the long-term goals of the research discipline. An ideal information extraction system should be *domain independent* or at least portable to any domain with a minimum amount of engineering effort. Moreover, Riloff and Lorenzen do not specify further the types of information. Although many different types of semantics can be defined, the semantics – whether they are defined in a specific or a general subject domain – ideally should be as much as possible *universally accepted* and *bear on the ontological nature and relationships of being*.

Another consequence of the stress on pattern matching approaches that were developed during the MUC competitions is that eventually any technique in which pattern matching is used to organize data in some structured format can be considered to be information extraction. For instance, the early nineties saw a sudden surge in popularity of research into approaches that try to extract the content of websites (e.g., shopbots that extract and compare prices), usually in order to convert them into a more convenient, uniform structural format. Some of these approaches analyze the natural language content of full text websites, but many only use pattern matching techniques that exploit the structural properties of markup languages to harvest data from automatically generated web pages. While many researchers conveniently gathered these approaches under the common denominator *web based information extraction* (see for instance Eikvil, 1999), we will assume that information extraction presupposes at least some degree of semantic content analysis. In addition, information extraction is also very much involved in finding the relationships that exist between the extracted information, based on evidence in text (e.g., John **kisses** Claudia).

Cowie and Lehnert try to mend the previous inaccuracies. They see information extraction as a process that involves the extraction of fragments of information from natural language texts and the linking of these fragments into a coherent framework. In their view, information extraction

[...] isolates relevant text fragments, extracts relevant information from the fragments, and then pieces together the targeted information in a coherent framework. [...] The goal of information extraction research is to build systems that find and link relevant information while ignoring extraneous and irrelevant information.

(Cowie and Lehnert, 1996, p. 81)

Cowie and Lehnert's interpretation of information extraction is close to what we need to solve the problems at the beginning of this chapter. There is still one thing missing in their definition. Although in this book we concentrate on information extraction from text, text is not the only source of unstructured information. Among these sources, it is probably the source where one has made the largest advancements in automatic understanding. But, other sources (e.g., image, video) exhibit a similar need for *semantically labeling unstructured information*, and advances in their automatic understanding are expected to occur in the near future. Any framework in which information extraction functions should not exclude this given.

The interpretations above are only a few representative definitions, and in the literature one finds additional variant definitions. To this multitude, we will add our own working definition, trying to incorporate the kernel task and function of information extraction and to avoid both Riloff and Lorenzen's and Cowie and Lehnert's limitations:

DEFINITION

Information extraction is the identification, and consequent or concurrent classification and structuring into semantic classes, of specific information found in unstructured data sources, such as natural language text, making the information more suitable for information processing tasks.

This definition is concise and covers exactly in what sense we will use the term *information extraction* throughout this book, but it is still fairly abstract. In the next sections, we will clarify what its constituent parts exactly mean.

1.2 Explaining Information Extraction

1.2.1 Unstructured Data

Information extraction is used to get some information out of *unstructured data*. Written and spoken text, pictures, video and audio are all forms of unstructured data. *Unstructured* does not imply that the data is structurally incoherent (in that case it would simply be nonsense), but rather that its information is encoded in such a way that makes it difficult for computers to immediately interpret it. It would actually be more accurate to use the

terms *computationally opaque* vs. *computationally transparent* data. Information extraction is the process that adds meaning to unstructured, raw data, whether that is text, images, video or audio. Consequently, the data become structured or semi-structured and can be more easily processed by the computer (e.g., in information retrieval, data mining or summarization).

In this book, the unstructured data sources we are mainly concerned with are *written natural language texts*. The texts can be of different type or genre (e.g., news articles, scientific treaties, police reports). Unless stated otherwise, we will assume that these texts are rather well formed, i.e., that they are largely coherent and error free. This is far from always the case. Written data, especially in an electronic context, is notorious for being incoherent and full of grammatical and spelling errors that are drafted with (e.g., spam messages) or without purpose (e.g., instant messages, postings on informal news groups). Errors might also occur in the output of an automatic speech recognition system. The algorithms described in this book can all be applied to these deviant types of textual data, provided that the system is specifically trained to deal with the characteristics of the relevant text type or that the probabilities of the translation into elements of well formed text are taken into consideration.

Limiting this book to information extraction from the text medium is no restriction of its value. The technologies described in this book contribute to an advanced understanding of textual sources, which, for instance can be used for aligning text and images when training systems that understand images. In addition, because technologies for the understanding of media other than text will be further developed in the near future, it seems valuable to compile information extraction technologies for text as they serve as a source of ideas for content recognition in other media or in combined media (e.g., images and text, or video).

1.2.2 Extraction of Semantic Information

Information extraction identifies information in texts by taking advantage of their linguistic organization. Any text in any language consists of a *complex layering* of recurring patterns that form a coherent, meaningful whole. This is a consequence of the *principle of compositionality* (Szabó, 2004), a general notion from linguistic philosophy that underlies many modern approaches to language and that states that the meaning of any complex linguistic expression is a function of the meanings of its constituent parts. An English sentence typically contains a number of constituent

parts (e.g., a subject, a verb, maybe one or more objects). Their individual meanings, ordering and realization (for instance, the use of a specific verb tense) allow us to determine what the sentence means. If a text would be completely irregular, it would simply be impossible for humans to make any sense of it.

It is yet not entirely clear how these linguistic layers exactly interact, but many linguistic theories and natural language processing assume the existence of a *realizational chain*. This theoretical notion has its roots in the grammar that was written by the Indian grammarian *Panini* in the 6th - 5th century B.C. (see Kiparsky, 2002). According to this notion, meaning in a language is realized in the linguistic surface structure through a number of distinct linguistic levels, each of which is the result of a projection of the properties of higher, more abstract levels. For instance, for Panini the meaning of a simple sentence starts as an idea in the mind of a writer. It then passes through the stage in which the event and all its participants are translated into a set of semantic concepts, each of which is in its turn translated in a set of grammatical and lexical concepts. These are in their turn translated into the character sequences that we see written down on a page of paper. Information extraction (and natural language processing, for that matter) assumes that this projection process is to a considerable extent bidirectional, i.e., that ideas are recoverable from their surface realizations by a series of inverse processes.

In other words, information extraction presupposes that although the semantic information in a text and its linguistic organization is not *immediately* computationally transparent, it can nevertheless be retrieved by taking into account surface regularities that reflect its computationally opaque internal organization. An information extraction system will use a set of extraction patterns, which are either manually constructed or automatically learned, to take information out of a text and put it in a more structured format. The exact techniques that are used to extract semantic information from a natural language text form the main topic of this book. Particular methodologies and algorithms will be discussed throughout its main chapters.

The use of the term *extraction* implies that the semantic target information is *explicitly present* in a text's linguistic organization, i.e., that it is readily available in the lexical elements (words and word groups), the grammatical constructions (phrases, sentences, temporal expressions, etc.) and the pragmatic ordering and rhetorical structure (paragraphs, chapters, etc.) of the source text. In this sense, information extraction is different from techniques that *infer* information from texts, for instance by building logical rules (logical inference) and by trying to distil world or domain

knowledge from the propositions in a text through deductive, inductive or abductive reasoning. We will refer to this latter kind of information as *knowledge*. Knowledge discovery is also possible by means of statistical *data mining* techniques that operate on the information extracted from the texts (also referred to as *text mining*). In all these operations information extraction is often an indispensable preprocessing step. For instance, information that is extracted from police reports could be used as the input for a data mining algorithm for profiling or for detecting general crime trends, or as the input of a case based reasoning algorithm that predicts the location of the next strike of a serial killer based on similar case patterns.

1.2.3 Extraction of Specific Information

Information extraction is traditionally applied in situations where it is *known in advance* which kind of semantic information is to be extracted from a text. For instance, it might be necessary to identify what kind of events are expressed in a certain text and at what moment these events take place. Since in a specific language, events and temporal expressions can only be expressed in a limited number of ways, it is possible to design a method to identify specific events and corresponding temporal location in a text. Depending on the information need, different models can be constructed to distinguish different kinds of classes at different levels of semantic granularity. In some applications, for example, it will suffice to indicate that a part of a sentence is a temporal expression, while in others it might be necessary to distinguish between different temporal classes, for instance between expressions indicating past, present and future.

Information extraction does not present the user with entire documents, but it extracts *textual units* or elements from the documents, typically simple or multi-term basic phrases (Appelt and Israel, 1999), which we also call *text regions*. As such, information extraction is different from *extractive summarization*, which usually retrieves entire sentences from texts that serve as its summary. Information extraction, however, can be a useful first step in *extractive headline summarization*, in which the summary sentence is further reduced to a string of relevant phrases similar to a newspaper headline.

Specificity implies that not only the semantic nature of the target information is predefined in an information extraction system, but also the *unit* and *scope* of the elements to be extracted. Typical *extraction units* for an extraction system are word compounds and basic noun phrases, but in some applications it might be opportune to extract other linguistic units,

such as verb phrases, temporal markers, clauses, strings of related meanings that persist throughout different sentences, larger rhetorical structures, etc. Whereas the unit of extraction has to do with the granularity of individual information chunks that are lifted out of the source text, the *scope* of extraction refers to the granularity of the extraction space for each individual information request. Information can be extracted from one clause or from multiple clauses or sentences spanning one or more texts before it is outputted by the system. Consider the example that an information question wants to retrieve event information about assassinations, it might be that the name of the person assassinated and the time and place of the event is named in a first sentence of a news article, but that the name of the assassin and his method are mentioned in some sentences further in the discourse.

During the Message Understanding Conferences (MUC), there gradually arose a set of typical information extraction tasks (see Grishman and Sundheim, 1996; Cunningham, 1997). A most popular task probably is *named entity recognition*, i.e., recognizing person names, organizations, locations, date, time, money and percents. These names are often expanded to protein names, product brands, etc. Other tasks are *event extraction*, i.e., recognizing events, their participants and settings, and *scenario extraction*, i.e., linking of individual events in a story line. *Coreference resolution*, i.e., determining whether two expressions in natural language refer to the same entity, person, time, place, and event in the world, also receives quite a lot of attention. These task definitions have been extremely influential in concurrent information extraction research and we will see that although they are getting too narrow to cover everything that is presently expected from information extraction, they still define its main targets. Currently, we see a lot of interest for the task of *entity relation recognition*. A number of *domain specific extractions* are also popular, e.g., extraction of the date of availability of a product from a Web page, extraction of scientific data from publications, and extraction of the symptoms and treatments of a disease from patient reports. The interest in the above extraction tasks is also demonstrated in the current *Automatic Content Extraction (ACE)* project.

1.2.4 Classification and Structuring

Typical for information extraction is that information is not just extracted from a text but afterwards also *semantically classified* in order to ensure its future use in information systems. By doing this, the information from unstructured text sources also becomes structured (i.e., computationally

transparent and semantically well defined). In the extreme case, the information that is verbatim extracted from the texts is discarded for further processing, but this is not what is usually intended.

Any classification process requires a semantic classification scheme, i.e., a set of semantic classes that are organized in some relevant way (for instance in a hierarchy) and that are used to categorize the extracted chunks of information into a number of meaningful groups. A very large variety of semantic classification schemes are conceivable, ranging from a small set of abstract semantic classes to a very elaborate and specific classification.

Based on the general information focus of a system, we can make a main distinction between *closed domain* and *open domain* (or *domain independent*) information extraction systems. Traditionally, information extraction systems were closed domain systems, which means that they were designed to function in a rather specialized, well delineated knowledge domain (and that they will therefore use very specific classification rules). For instance, most MUC systems covered very limited subjects such as military encounters, Latin-American terrorism or international joint ventures (Grishman and Sundheim, 1996). Domain independent information extraction systems, on the other hand, are capable of handling texts belonging to heterogeneous text types and subject domains, and usually use very generic classification schemes, which might be refined, if the information processing task demands a more specific identification of semantic information. The technology described in this book applies to both closed and open domain information extraction.

We mentioned before that information extraction essentially converts unstructured information from natural language texts into structured information. This implies that there has to be a predefined structure, a representation, in which the extracted information can be cast. Although the extracted information can solely be labeled for consequent processing by the information system, in the past many template based extraction systems have been developed. Template representations were typically used to describe single events (and later also complex scenarios) and consist of a set of attribute-value pairs (so-called *slots*), each of which represents a relevant aspect of the event (e.g., the action or state, the persons participating, time, place). An information extraction task traditionally tries to take information from a source text and maps it to an empty value of the defined template.

In order to know which piece of information is supposed to end up in which template slot, an information extraction application uses a set of extraction rules. These rules state which formal or linguistic properties a particular chunk of information must possess to belong to a particular

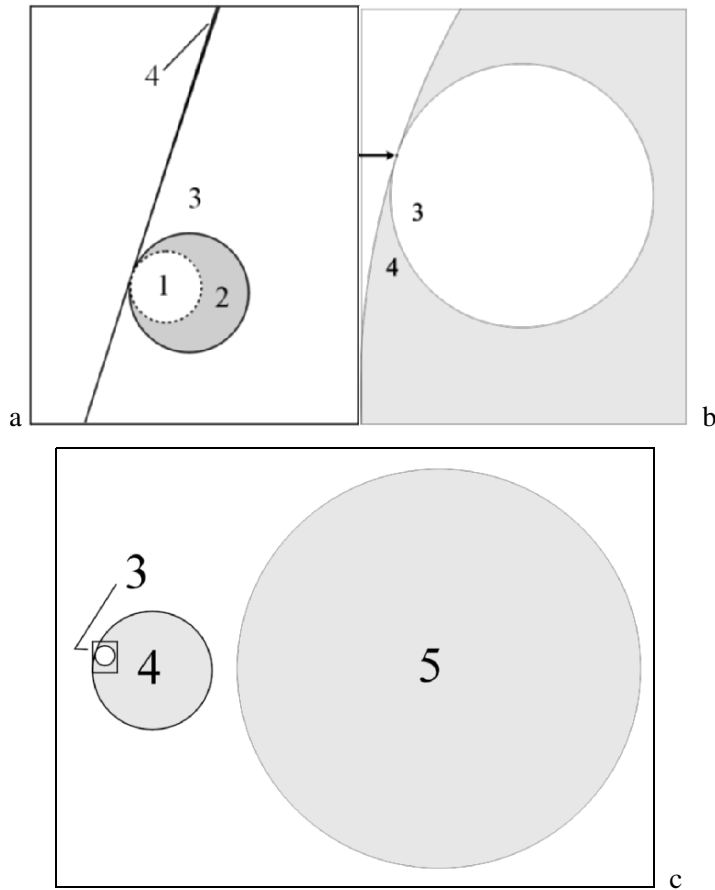
semantic class. Especially in earlier systems these rules were usually handcrafted (e.g., the FASTUS system developed by Appelt et al., 1993). Currently, machine learning is playing a central role in the information extraction paradigm. In most cases, *supervised learning* is used, in which a learning algorithm uses a training corpus with manually labeled examples to induce extraction rules, as they are applicable to a particular language and text type (e.g., the CRYSTAL system developed by Soderland et al. 1995). In some cases, it is also possible to apply *unsupervised learning*, for which no training corpus is necessary. For instance, unsupervised learning systems have been implemented for noun phrase coreferent resolution (e.g., Cardie and Wagstaff, 1999). Today, we see a large interest in weakly supervised learning approaches that limit the number of examples to be manually labeled (Shen et al., 2004). The application of these learning techniques has been one of the main enabling factors for information extraction to move from very domain specific to more domain independent analyses. In addition, the machine learning techniques more easily allow modeling a probabilistic class assignment instead of a purely deterministic one.

1.3 Information Extraction and Information Retrieval

1.3.1 Information Overload

Our modern world is flooded with information. Nobody exactly knows how much information there is – or how one could uniformly measure information flows from heterogeneous sources – but Lyman and Hal (2003) estimate that the total amount of newly created information on physical media (print, film, optical and magnetic storage) amounted to some 5 exabytes in 2002, most of it is stored in digital format. This corresponds to 9,500 billion books or 500,000 times the entire Library of Congress (which is supposed to contain approximately 10 terabytes of information). According to their measures, the surface Web contains around 167 terabytes of information, and there are indications that the deep Web, i.e., information stored in databases that is accessible to human users through query interfaces but largely inaccessible to automatic indexing, is about 400 to 500 times larger. Lyman and Hal (2003) estimate it to be at least 66,800 terabytes of data. A large fraction of this information is unstructured in the form of text, images, video and audio. These gargantuan figures are already

outdates at the time of writing and are dwarfed by the amount of e-mail traffic that is generated, which according to Lyman and Hal (2003) amounts to more than 300,000 terabytes of unique information per year.



Legend

- | | |
|---|--|
| 1 | Size of the Library of Congress |
| 2 | Size of the surface Web |
| 3 | Size of the surface + deep Web |
| 4 | Size of surface + deep Web + e-mail traffic |
| 5 | Size of text data on hard discs sold in 2003 |

Fig. 1.1. Graphical presentation of the size of the Web and global storage capacity on computer hard discs anno 2003.

According to these authors, during 2003 an estimated 15,892.24 exabytes of hard disc storage was sold worldwide. A similar study that confirms the information overload is made by O'Neill, Lavoie and Bennett (2003). If this trend continues we will have to express amounts of information in yottabytes (2^{80} number of bytes).

In order to give a rough impression about the amounts of data that are involved, Fig. 1.1 gives a graphical representation of the total amount of data present on the Web and on hard discs worldwide in 2003. Data ratios are reflected in the relative size differences between the diameters of the circles. Figure 1.1a shows the size of unique textual data on the surface web [2] in comparison with the textual data on the combined surface and deep web [3] and of the surface and deep Web plus all e-mail traffic [4]. The size of the Library of Congress [1] is given as a reference. Figure 1.1b is a 180-fold magnification in which Fig. 1.1a appears as the minute rectangle at the point of tangency of 3 and 4. Fig. 1.1c gives an impression of the complete size of the textual web at the left hand side and a comparison with all textual data on hard discs on the right hand side.

This immense information production prevents its users to efficiently select information and accurately use it in problem solving and decision making (Edmunds and Morris, 2000; Farhoomand and Drury, 2002). Even if we find ways for reducing the information generation, there still is a large demand for intelligent tools that assist humans in information selection and processing (Berghel, 1997).

1.3.2 Information Retrieval

Information retrieval (IR) is a solution to this kind of problem (Baeza-Yates and Ribeiro-Neto, 1999). It allows a user to retrieve a set of documents from large document collections, such as the Web or a corporate intranet, based on a keyword based query. Information retrieval is able to search efficiently through huge amounts of data because it builds indexes from the documents in advance in order to reduce the time complexity of each real-time search. The low level keyword matching techniques that are generally used in information retrieval systems make them error tolerant, domain independent and – above all – very fast (Lewis and Sparck Jones, 1996). The success of information retrieval systems in general, and the Web search engines in particular is largely due to the flexibility of these systems with regard to the queries that users pose. Users have all kinds of information needs that are very difficult to determine a priori. Because users do not always pose their queries with the words that occur in relevant documents, query expansion with synonym and related terms is very

popular, primarily enhancing the recall of the results of the search. Information retrieval is very successful in what it is aimed to do, namely providing a rough and quick approach to find relevant documents.

A downside is that such a robust and flexible approach sometimes results in a low precision of an information search and in a huge number of possibly relevant documents when a large document base is searched, which are impossible to consult by the user of the information system (Blair, 2002).

1.3.3 Searching for the Needle

Because of the information overload, the classical information retrieval paradigm is no longer preferable. This *paradigm* has found its roots in the example of the *traditional library*. One is helped in finding potentially relevant books and documents, but the books and documents are still consulted by the humans. When the library is becoming very large and the pile of potentially relevant books is immensely high, humans want more advanced information technology to assist them in their information search. We think that information extraction technology plays an important role in such a development.

Currently, an information retrieval system returns a list of relevant documents, where each individual document has to be fetched and skimmed through in order to assess its real relevance. There is a need for tools that reduce the amount of text that has to be read to obtain the desired information. To address this need, the information retrieval community is currently exploring ways of pinpointing highly relevant information. This is one of the reasons *question answering systems* are being researched. The user of a question answering retrieval system expresses his or her information need as a natural language question and the system extracts the answer to the information question from the texts of the documents (Maybury, 2003).

Information extraction is one of the core technologies to help facilitate highly focused retrieval. Indeed, recognizing entities and semantically meaningful relations between those entities is a key to provide focused information access.

With the current interest in expressing queries as natural language texts, the need for semantic classification of entities and their relations in the texts of document and query becomes of primordial importance. Information extraction technology realizes that – simply saying – not only the words of the query, but also the semantic classifications of entities and

their relations must match the information found in the documents (Moens, 2002).

Especially in information gathering settings where the economical costs of searching is high or in time critical applications such as military or corporate intelligence gathering, a user often needs very specific information very quickly. For instance, an organization might need a list of all companies that have offices in the Middle East and conducted business transactions or pre-contract negotiations in the Philippines or Indonesia in the last five months. He or she knows that many of this information is available in news feeds that were gathered over the last half year, but it is impossible to go through tens of thousands of news snippets to puzzle all relevant data together. In addition, we cannot neglect the need for flexible querying. There will always be a large variety of dynamically changing information needs.

Information retrieval techniques typically use general models for processing large volumes of text. The *indices* are stored in data structures that are especially designed to be efficiently searched at the time of querying. An ideal information retrieval system answers all kinds of possible information questions in a very precise way by extracting the right information from a (possibly large) collection of documents.

Information extraction helps building such information systems. The extracted information is useful to construct sensitive indices more closely linked to the actual meaning of a particular text (Cowie and Lehnert, 1996). This is often only restricted to the recognition and classification of entities that are referenced in different places of the text and recognition of relations between them. Besides an index of words that occur in the documents, certain words or other information units are tagged with additional semantic information. This meta-information allows more precisely answering information questions without losing the advantages of flexible querying. Information extraction technology allows for a much richer indexing representation of both query and document or information found in the document, which can improve retrieval performance of both open and closed domain texts. Especially linguistically motivated categories of semantics become important (e.g., expressions of time, location, coreference, abstract processes and their participants, ...). As we will show in this book, the identified and classified information – even if very generic semantic classifications are made – is useful in information retrieval and selection, allowing for answers of information needs that can be more precisely inferred from information contained in documents. Information extraction can be regarded as

a kind of cheap and easy form of natural language understanding, which can be integrated in an information retrieval system to roughly provide some understanding of query and document.

As such, information extraction introduces *natural language processing* (NLP) technology into retrieval systems. Natural language processing has sought models for analyzing human language. These attempts were sometimes successful, but they also aroused an awareness of the enormous magnitude of the task. NLP deals with the processing of the linguistic structure of text. This includes morphological, syntactic and semantic analysis of language, the extraction of discourse properties and domain knowledge or world knowledge, and eventually aims at full natural language understanding. The problems of automatic text understanding concern the encoding of the necessary knowledge into a system and constructing a proper inference mechanism, as well as the computational complexity of the necessary operations. Information retrieval technology traditionally has depended upon relatively simple, but robust methods, while natural language processing involves complex knowledge based systems that have never approached robustness. The much more recent field of *information extraction*, which is a first step towards full natural language understanding, reaches a degree of maturity and robustness that makes it ready to incorporate in information processing systems.

Also in *Cross-Language Information Retrieval* (CLIR), information extraction plays an important role. In CLIR, a query in one language searches a document base in another language. The semantic concepts that information extraction uses are mostly language independent, but more accurately translate a query and map the translated queries with the documents. To minimize the language specific alterations that need to be made in extending an information extraction system to a new language, it is important to separate the task specific conceptual knowledge the system uses, which may be assumed to be language independent, from the language dependent lexical knowledge the system requires, which unavoidably must be extended or learned for each new language. The language independent domain model can be compared to the use of an interlingua representation in Machine Translation (MT). An information extraction system however does not require full generation capabilities from the intermediate representation (unless the extracted information is translated) and the task will be well specified by a limited model of semantic classes, rather than a full unrestricted world model. This makes an interlingua representation feasible for information extraction.

1.4 Information Extraction and Other Information Processing Tasks

People are not only interested in information retrieval systems that precisely give an answer to an information query, they also use information systems that help in solving problems, such as data mining systems, systems that reason with knowledge, and systems that visualize, synthesize, or summarize information. Given the current information overload, system assistance is more than welcome. Information extraction is a helpful step in these processes because the data become structured and semantically enriched.

A typical example is applying data mining techniques to the information found in texts. Examples are law texts and police reports that are parsed in order to analyze specific trends (Zeleznikow and Stranieri, 2005). There is an increasing interest in extracting knowledge from texts to be used in knowledge-based systems. A knowledge based system uses knowledge that is formally represented in a knowledge representation language and reasons with this knowledge in the search to find an answer to an information question. The knowledge can be in the form of sharable knowledge components such as an ontology or in the form of very specific knowledge that is used for performing a specific task. Information extraction technology is very useful to automatically build the knowledge rules and frames. We can refer to an example of the legal domain. A nice example is the automatic translation of legislation into knowledge rules. The rules can be used to infer the answer to a specific problem (e.g., **Is the cultivation of Erythroxylon punishable?**) (Moens, 2003). Similarly, technical documentation can be automatically translated into knowledge structures to be questioned at the time a problem occurs.

Information extraction *semantically classifies* textual units. As such information extraction is related to *text categorization* and *abstractive summarization*. Text categorization classifies text passages or complete texts with semantic labels. These labels are usually used in the matching of an information query with the document in a retrieval context or for filtering documents according to a certain user profile. Abstractive summarization will replace a text or a text passage by a more abstract concept or concepts. Although text categorization, abstractive summarization and information extraction techniques overlap to a large degree, the semantic labels in text categorization and abstracting define the most salient information of a text in one or a few abstract concepts. Information extraction is here somewhat the opposite as it allows finding detailed information, for instance, with regard to a certain event. On the other hand information extraction and text

categorization complement each other in two directions. In a top down approach, very domain specific information extraction technologies can be selected based on a prior semantic classification of a complete text or text passage. In a bottom up approach the detailed information labeled with information extraction technologies can contribute to a more fine-grained classification of a complete text or passage.

There is a large need to synthesize and summarize information that is, for instance, collected as the result of an information search. Information extraction technologies identify the entities that are involved in certain processes and the relations between entities. Such information better permits generating concise headlines or compressed sentences that make up the summary.

1.5 The Aims of the Book

The main goal of the book is to give a *comprehensive overview of algorithms used in information extraction from texts*. Almost equal importance is given to early technologies developed in the field primarily with the aim of natural language understanding, as to the most advanced and recent technologies for information extraction. The past approaches are an incentive to identify some forgotten avenues that can be researched to advance the state of the art. Machine learning is playing a central role in the development of the novel technologies and contributes to the portability and widespread use of information extraction technology. The book will especially focus on weakly supervised learning algorithms that are very promising for information extraction.

A second important aim is to focus on the *prospects* of information extraction *to be used* in modern information systems, and more specifically in *information retrieval systems*. We want to demonstrate that on one hand the statistical and machine learning techniques traditionally used in information retrieval have little attention for the underlying cognitive and linguistic models that shape the patterns that we are attempting to detect. On the other hand, current models of information retrieval that use expressive query statements in natural language (e.g., simple and complex question answering, a textual query by example) exhibit the need for a semantic based matching between the content of the query and the documents.

It is also argued that information extraction technology has evolved from a template extraction technique to a real *necessary labeling step* in many *information management* tasks. In the past, as a result of the development of domain specific information extraction technology lexicons

were built that store subcategorization patterns for domain specific verbs in such a fashion as to permit them to use the patterns directly in a pattern-matching engine. Here we look to information extraction differently. We see information extraction as a tool that aids in other tasks such as information retrieval. In such a framework, information extraction is seen as a kind of preprocessing and labeling of texts, which contributes to the performance of other tasks (here mainly information retrieval, but we also refer to data mining, knowledge discovery, and summarization). We aim at labeling information in open domain and closed domain texts, at identifying specific facts or more generic semantic information, and store this information in a format that allows flexible further processing. This evolution is sustained by content recognition techniques that are currently developed for other media such as images for which the information extraction paradigm also applies.

Whereas traditional information extraction recognizes information in text in a deterministic way and represents this information in a format with known semantics, such as relational database fields, we leave room for *probabilistic classification of the information* and consequent probabilistic processing of the information (e.g., in a probabilistic retrieval model). This is an approach that better corresponds with the way we humans search for information in texts of a language we do not completely understand, but from the combined combination of evidence we can make a good guess and fairly accurately locate the information. This is also an approach that better fits the philosophy and tradition of information retrieval.

In addition, we want to demonstrate that the form of semantic labeling that is advocated by information extraction technologies does not put in danger the cherished *flexibility of an information search*. According to our definition that is given p. 4, the extracted information is classified and structured. The information retrieval models can smoothly integrate this structured information when matching the information question with the document. In this respect current *XML* (Extensible Markup Language) retrieval models that combine and rank structured and unstructured information are a source of inspiration (Blanken et al., 2003).

Last but not least, we want to illustrate the information extraction technologies and evaluate the results as they are currently implemented in many different application domains. The illustrative examples intend to demonstrate the generic character of current information extraction technologies and their portability to different domains. The evaluation should also lead to better insights into the points of attention for future research.

We focus in the book on information extraction from text in natural language. Although information extraction takes into account certain characteristics of the language and the examples are usually taken from the English

language, our discussion is as language independent as possible. Many text documents currently have structure and layout tagged in markup languages such as XML and HTML (HyperText Markup Language). Such markups are not the focus of our attention.

The book is organized as follows.

In *Chapter 2* we give a short historical overview of information extraction, starting from its early origins in the mid seventies that coincide with early attempts of natural language understanding, where we do not ignore the influence of artificial intelligence authorities such as Roger Schank and Marvin Minsky. During this period the use of symbolic, handcrafted knowledge was very popular. We explain the general trend towards machine learning techniques that started in the early nineties and finish with the current interest in weakly supervised learning methods and hybrid technologies that combine the use of symbolic and machine learning techniques. The dive into a history of more than three decades focuses on the different factors of success and causes of difficulties. This chapter also explains typical information extraction tasks from both a linguistic theoretical and application oriented viewpoint. The tasks among others include named entity recognition, noun phrase coreferent resolution, semantic role classification, and the recognition and resolution of temporal expressions. The chapter also describes the general architecture of an information extraction system and assisting linguistic and knowledge resources.

Chapter 3 gives an in depth discussion of some of the most important symbolic techniques for information extraction that use handcrafted knowledge. We start from the Conceptual Dependency Theory of Roger Schank, explain in detail frame based approaches and the use of finite state automata to parse the texts.

Chapter 4 offers an introduction to the current pattern recognition methods that use machine learning methods. A substantial part of this chapter is devoted to the features of the texts that are used in the information extraction tasks. The features include lexical, syntactic, semantic and discourse features found in the texts as well as features derived from external knowledge resources.

Chapter 5 explains the most important and most successful supervised machine learning techniques currently in use in information extraction. We explain Support Vector Machines, maximum entropy modeling, hidden Markov models, conditional random fields, learning of decision rules and trees, and relational learning. The theoretical background of a technology is illustrated with realistic examples of information extraction tasks.

Chapter 6 is devoted to unsupervised learning aids. Such aids become very popular because of the high cost of manual annotation. The techniques described range from completely unsupervised methods such as

clustering to weakly supervised methods such as co-training, self-training and active learning. Again the theory of each technique is illustrated with a real information extraction example.

Chapter 7 integrates information extraction in an information retrieval framework. More specifically, it studies how information extraction is incorporated in the various existing retrieval models. A retrieval model matches query and document representations and ranks the documents or information found in the documents according to relevance to the query. The different models that are discussed are the classical vector space model, the language model, the inference network model and the logic based model. Because information extraction leaves behind a *bag-of-words* representation of query and documents, we need adapted indexing structures that can be efficiently searched at the time of querying.

Chapter 8 discusses the evaluation metrics currently in use in information extraction. Evaluation metrics allow a comparison of technologies and systems. Classical metrics such as recall, precision and accuracy are discussed. Also metrics that value the goodness of a clustering are important to mention (e.g., to evaluate noun phrase coreferent resolution). We mention also the evaluation metrics currently in use in international competitions. The difference between an intrinsic and an extrinsic evaluation of information extraction is explained.

Chapter 9 elaborates on many recent applications in information extraction and gives the reader a good assessment of the capabilities of current technologies. Information extraction is illustrated with applications of news services, intelligence gathering and extracting content from the biomedical, business and legal domains. Finally, we study the special case of information extraction from noisy texts (e.g., transcribed speech) and its difficulties.

Chapter 10 summarizes our most important findings with regard to the algorithms used in information extraction and the future prospects of this technology. A section will also elaborate on future promising improvements of the technologies.

Each chapter is accompanied by clear and illustrative examples, and the most relevant past and current bibliography is cited.

1.6 Conclusions

In this first chapter we have defined information extraction from text as a technology that identifies, structures and semantically classifies certain information in texts. We have demonstrated the importance of information

extraction in many information processing tasks, among which information retrieval. In the next chapter, we give an historical overview of the information extraction technologies, define in detail typical information extraction tasks and outline the architecture of an information extraction system. This will give the reader a better understanding of information extraction needs as they have arisen in the course of the last decades and will smoothly introduce the different technologies that are discussed in the main parts of the book.

1.7 Bibliography

- ACE: www.nist.gov/speech/tests/ace/
- Appelt, Douglas E. and David J. Israel (1999). Introduction to information extraction technology. *Tutorial at the International Joint Conference on Artificial Intelligence IJCAI-99*: <http://www.ai.sri.com/~appelt/ie-tutorial/>
- Appelt, Douglas E., Jerry R. Hobbs, John Bear, David J. Israel and Mabry Tyson (1993). FASTUS: A finite-state processor for information extraction from real-world text. In *Proceedings of the Thirteenth International Joint Conference on Artificial Intelligence* (pp. 1172-1178). San Mateo, CA: Morgan Kaufmann.
- Baeza-Yates, Ricardo and Berthier Ribeiro-Neto (1999). *Modern Information Retrieval*. Harlow, UK: Addison-Wesley.
- Berghel, Hal (1997). Cyberspace 2000: Dealing with information overload. *Communications of the ACM*, 40 (2), 19-24.
- Blair, David C. (2002). The challenge of commercial document retrieval, part I: Major issues, and a framework based on search exhaustivity, determinacy of representation and document collection size. *Information Processing and Management*, 38 (2), 273-291.
- Blanken, Henk M., Torsten Grabs, Hans-Jörg Schek, Ralf Schenkel and Gerhard Weikum (Eds.) (2003). *Intelligent Search on XML Data, Applications, Languages, Models, Implementations and Benchmarks (Lecture Notes in Computer Science, 2818)*. New York, NY: Springer.
- Cardie, Claire and Kiri Wagstaff (1999). Noun phrase coreference as clustering. In *Proceedings of the Joint Conference on Empirical Methods in Natural Language Processing and Very Large Corpora* (pp. 82-89). East Stroudsburg, PA: ACL.
- Cowie, Jim and Wendy Lehnert (1996). Information extraction. *Communications of the ACM*, 39 (1), 80-91.
- Cunningham, Hamish (1997). *Information Extraction: A User Guide*. Research memo CS-97-02. Sheffield: University of Sheffield, ILASH.
- Edmunds, Angela and Anne Morris (2000). The problem of information overload in business organisations: A review of the literature. *International Journal of Information Management*, 20, 17-28.
- Eikvil, Line (1999). *Information Extraction from the World Wide Web: A Survey*. Norwegian Computer Center, Report no. 945, July 1999.

- Farhoomand, Ali F. and Don H. Drury (2002). Managerial information overload. *Communications of the ACM*, 45 (10), 127-131.
- Grishman, Ralph and Beth Sundheim (1996). Message Understanding Conference 6: A brief history. In *Proceedings of the 16th International Conference on Computational Linguistics* (pp. 466-471). San Mateo, CA: Morgan Kaufmann.
- Kiparsky, Paul (2002). *On the Architecture of Panini's Grammar*. Three lectures delivered at the Hyderabad Conference on the Architecture of Grammar, January 2002, and at UCLA, March 2002.
- Lewis, David D. and Karen Sparck Jones (1996). Natural language processing for information retrieval. *Communications of the ACM* 39 (1), 92-101.
- Lyman, Peter and Hal R. Varian (2003). *How Much Information? 2003*. URL: <http://www.sims.berkeley.edu/research/projects/how-much-info-2003/>
- Maybury, Mark (2003) (Ed.). *New Directions in Question Answering*. In *Papers from the 2003 AAAI Spring Symposium*. Menlo Park, CA: The AAAI Press.
- Moens, Marie-Francine (2002). What information retrieval can learn from case-based reasoning. In *Proceedings JURIX 2002: The Fifteenth Annual Conference (Frontiers in Artificial Intelligence and Applications)* (pp. 83-91). Amsterdam: IOS Press.
- Moens, Marie-Francine (2003). Interrogating legal documents: The future of legal information systems? In *Proceedings of the JURIX 2003 Workshop on Question Answering for Interrogating Legal Documents December 11, 2003* (pp. 19-30). Utrecht University, The Netherlands.
- O'Neill, Edward T., Brian F. Lavoie and Rick Bennett (2003). Trends in the evolution of the public web. 1998-2002. *D-Lib Magazine* 9 (4).
- Riloff, Ellen and Jeffrey Lorenzen (1999). Extraction-based text categorization: Generating domain-specific role relationships automatically. In Tomek Strzalkowski (Ed.), *Natural Language Information Retrieval* (pp. 167-196). Dordrecht, The Netherlands: Kluwer Academic Publishers.
- Shen, Dan, Jie Zhang, Jian Su, Guodong Zhou and Chew-Lim Tan (2004). Multi-criteria-based active learning for named entity recognition. In *Proceedings of the 42nd Annual Meeting of the Association for Computational Linguistics* (pp. 590- 597). East Stroudsburg, PA: ACL.
- Soderland, Stephen, David Fisher, Jonathan Aseltine and Wendy Lehnert (1995). CRYSTAL: Inducing a conceptual dictionary. In *Proceedings of the Fourteenth International Joint Conference on Artificial Intelligence* (pp. 1314-1319). San Mateo, CA: Morgan Kaufmann.
- Sundheim, Beth M. (1992). Overview of the fourth Message Understanding evaluation and Conference. In *Proceedings of the Fourth Message Understanding Conference (MUC-4)* (pp. 3-21). San Mateo: CA: Morgan Kaufmann.
- Szabó, Zoltán Gendler (2004). Compositionality. In Edward N. Zalta (Ed.), *The Stanford Encyclopedia of Philosophy (Fall 2004 Edition)*.
- Zeleznikow, John and Andrew Stranieri (2005). *Knowledge Discovery from Legal Databases*. New York, NY: Springer.