

CHAPTER 1

INTRODUCTION

CHAPTER OVERVIEW

This chapter gives an idea of natural language processing (NLP) and information retrieval (IR). Various levels of analysis involved in NLP along with the knowledge used by these levels of analysis are discussed. Some of the difficulties in analysing text and specific factors that make automatic processing of languages difficult are also touched upon. The chapter underlines the role of grammar in language processing and introduces transformational grammar. Indian languages differ a lot from English. These differences are clearly pointed out. Further, a number of NLP applications are introduced along with some of the early NLP systems. Towards the end, information retrieval is discussed.

1.1 WHAT IS NATURAL LANGUAGE PROCESSING (NLP)

Language is the primary means of communication used by humans. It is the tool we use to express the greater part of our ideas and emotions. It shapes thought, has a structure, and carries meaning. Learning new concepts and expressing ideas through them is so natural that we hardly realize how we process natural language. But there must be some kind of representation in our mind, of the content of language. When we want to express a thought, this content helps represent language in real time. As children, we never learn a computational model of language, yet this is the first step in the automatic processing of languages. *Natural language processing* (NLP) is concerned with the development of computational models of aspects of human language processing. There are two main reasons for such development:

1. To develop automated tools for language processing
2. To gain a better understanding of human communication

Building computational models with human language-processing abilities requires a knowledge of how humans acquire, store, and process language. It also requires a knowledge of the world and of language.

Historically, there have been two major approaches to NLP—the *rationalist* approach and the *empiricist* approach. Early NLP research took a rationalist approach, which assumes the existence of some language faculty in the human brain. Supporters of this approach argue that it is not possible for children to learn a complex thing like natural language from limited sensory inputs. Empiricists do not believe in existence of a language faculty. Instead, they believe in the existence of some general organization principles such as pattern recognition, generalization, and association. Learning of detailed structures can, therefore, take place through the application of these principles on sensory inputs available to the child.

1.2 ORIGINS OF NLP

Natural language processing sometimes mistakenly termed natural language understanding—originated from machine translation research. While natural language understanding involves only the interpretation of language, natural language processing includes both understanding (interpretation) and generation (production). The NLP also includes speech processing. However, in this book, we are concerned with text processing only, covering work in the area of computational linguistics, and the tasks in which NLP has found useful application.

Computational linguistics is similar to theoretical- and psycho-linguistics, but uses different tools. Theoretical linguists mainly provide structural description of natural language and its semantics. They are not concerned with the actual processing of sentences or generation of sentences from structural description. They are in a quest for principles that remain common across languages and identify rules that capture linguistic generalization. For example, most languages have constructs like noun and verb phrases. Theoretical linguists identify rules that describe and restrict the structure of languages (grammar). Psycholinguists explain how humans produce and comprehend natural language. Unlike theoretical linguists, they are interested in the representation of linguistic structures as well as in the process by which these structures are produced. They rely primarily on empirical investigations to back up their theories. (Computational linguistics is concerned with the study of language using computational models of linguistic phenomena. It deals with the application of linguistic theories and computational techniques for NLP. In computational linguistics, representing a language is a major problem; most knowledge representations tackle only a small part of knowledge.

Representing the whole body of knowledge is almost impossible. The words knowledge and language should not be confused. This is discussed in detail in Section 1.3.

Computational models may be broadly classified under knowledge-driven and data-driven categories. Knowledge-driven systems rely on explicitly coded linguistic knowledge, often expressed as a set of handcrafted grammar rules. (Acquiring and encoding such knowledge is difficult and is the main bottleneck in the development of such systems.) They are, therefore, often constrained by the lack of sufficient coverage of domain knowledge. Data-driven approaches presume the existence of a large amount of data and usually employ some machine learning technique to learn syntactic patterns. The amount of human effort is less and the performance of these systems is dependent on the quantity of the data. These systems are usually adaptive to noisy data.

As mentioned earlier, this book is mainly concerned with computational linguistics approaches. We try to achieve a balance between semantic (knowledge-driven) and data-driven approaches on one hand, and between theory and practice on the other. It is at this point that the book differs significantly from other textbooks in this area. The tools and techniques have been covered to the extent that is needed to build sufficient understanding of the domain and to provide a base for application.

The NLP is no longer confined to classroom teaching and a few traditional applications. With the unprecedented amount of information now available on the web, NLP has become one of the leading techniques for processing and retrieving information. In order to cope with these developments, this book brings together information retrieval with NLP. The term information retrieval is used here in a broad manner to include a number of information processing applications such as information extraction, text summarization, question answering, and so forth. The distinction between these applications is made in terms of the level of detail or amount of information retrieved. We consider retrieval of information as part of processing. The word ‘information’ itself has a much broader sense. It includes multiple modes of information, including speech, images, and text. However, it is not possible to cover all these modes due to space constraints. Hence, this book focuses on textual information only.

1.3 LANGUAGE AND KNOWLEDGE

Language is the medium of expression in which knowledge is deciphered. We are not competent enough to define language and knowledge and its

implications. We are here considering the text from of the language and the content of it as knowledge.

Language, being a medium of expression, is the outer form of the content it expresses. The same content can be expressed in different languages. But can language be separated from its content? If so, how can the content itself be represented? Generally, the meaning of one language is written in the same language (but with a different set of words). It may also be written in some other, formal, language. Hence, to process a language means to process the content of it. As computers are not able to understand natural language, methods are developed to map its content in a formal language. Sometimes, formal language content may have to be expressed in a natural language as well. Thus, in this book, language is taken up as a knowledge representation tool that has historically represented the whole body of knowledge and that has been modified, maybe through generation of new words, to include new ideas and situations. The language and speech community, on the other hand, considers a language as a set of sounds that, through combinations, conveys meaning to a listener. However, we are concerned with representing and processing text only. Language (text) processing has different levels, each involving different types of knowledge. We now discuss various levels of processing and the types of knowledge it involves.

The simplest level of analysis is *lexical analysis*, which involves analysis of words. Words are the most fundamental unit (syntactic as well as semantic) of any natural language text. Word-level processing requires morphological knowledge, i.e., knowledge about the structure and formation of words from basic units (morphemes). The rules for forming words from morphemes are language specific.

The next level of analysis is *syntactic analysis*, which considers a sequence of words as a unit, usually a sentence, and finds its structure. Syntactic analysis decomposes a sentence into its constituents (or words) and identifies how they relate to each other. It captures grammaticality or non-grammaticality of sentences by looking at constraints like word order, number, and case agreement. This level of processing requires syntactic knowledge, i.e., knowledge about how words are combined to form larger units such as phrases and sentences, and what constraints are imposed on them. Not every sequence of words results in a sentence. For example, 'I went to the market' is a valid sentence whereas 'went the I market to' is not. Similarly, 'She is going to the market' is valid, but 'She are going to the market' is not. Thus, this level of analysis requires detailed knowledge about rules of grammar.

Yet another level of analysis is *semantic analysis*. Semantics is associated with the meaning of the language. Semantic analysis is concerned with creating meaningful representation of linguistic inputs. The general idea of semantic interpretation is to take natural language sentences or utterances and map them onto some representation of meaning. Defining meaning components is difficult as grammatically valid sentences can be meaningless. One of the famous examples is, 'Colorless green ideas sleep furiously' (Chomsky 1957). The sentence is well-formed, i.e., syntactically correct, but semantically anomalous. However, this does not mean that syntax has no role to play in meaning. Bach (2002) considers:

... semantics to be a projection of its syntax. That is semantic structure is interpreted syntactic structure.'

But definitely, syntax is not the only component to contribute meaning. Our conception of meaning is quite broad. We feel that humans apply all sorts of knowledge (i.e., lexical, syntactic, semantic, discourse, pragmatic, and world knowledge) to arrive at the meaning of a sentence. The starting point in semantic analysis, however, has been lexical semantics (meaning of words). A word can have a number of possible meanings associated with it. But in a given context, only one of these meanings participates. Finding out the correct meaning of a particular use of word is necessary to find meaning of larger units. However, the meaning of a sentence cannot be composed solely on the basis of the meaning of its words. Consider the following sentences:

Kabir and *Ayan* are married.
Kabir and *Suha* are married.

Both sentences have identical structures, and the meanings of individual words are clear. But most of us end up with two different interpretations. We may interpret the second sentence to mean that Kabir and Suha are married to each other, but this interpretation does not occur for the first sentence. Syntactic structure and compositional semantics fail to explain these interpretations. We make use of pragmatic information. This means that semantic analysis requires pragmatic knowledge besides semantic and syntactic knowledge.

A still higher level of analysis is *discourse analysis*. Discourse-level processing attempts to interpret the structure and meaning of even larger units, e.g., at the paragraph and document level, in terms of words, phrases, clusters, and sentences. It requires the resolution of anaphoric references and identification of discourse structure. It also requires discourse knowledge, that is, knowledge of how the meaning of a sentence is determined by preceding sentences—e.g., how a pronoun refers to the

preceding noun—and how to determine the function of a sentence in the text. In fact, pragmatic knowledge may be needed for resolving anaphoric references. For example, in the following sentences, resolving the anaphoric reference ‘they’ requires pragmatic knowledge:

- The district administration refused to give the trade union permission for the meeting because they feared violence.
 The district administration refused to give the trade union permission for the meeting because they oppose government.

The highest level of processing is *pragmatic analysis*, which deals with the purposeful use of sentences in situations. It requires knowledge of the world, i.e., knowledge that extends beyond the contents of the text. The Cyc project (Lenat 1986) at University of Austin is an attempt to utilize world knowledge in NLP. However, its usefulness in a general-domain NLP system is yet to be demonstrated. Furthermore, whether or not semantics can be associated with a symbol manipulator and whether humans use logic in the same way as the Cyc project, are both issues of debate.

1.4 THE CHALLENGES OF NLP

There are a number of factors that make NLP difficult. These relate to the problems of representation and interpretation. Language computing requires precise representation of content. Given that natural languages are highly ambiguous and vague, achieving such representation can be difficult. The inability to capture all the required knowledge is another source of difficulty. It is almost impossible to embody all sources of knowledge that humans use to process language. Even if this were done, it is not possible to write procedures that imitate language processing as done by humans. In this section, we detail some of the problems associated with NLP.

Perhaps the greatest source of difficulty in natural language is identifying its semantics. The principle of compositional semantics considers the meaning of a sentence to be a composition of the meaning of words appearing in it. In the earlier section, we saw a number of examples where this principle failed to work. Our viewpoint is that words alone do not make a sentence. Instead, it is the words as well as their syntactic and semantic relation that give meaning to a sentence. As pointed out by Wittgenstein (1953): ‘The meaning of a word is its use in the language.’ A language keeps on evolving. New words are added continually and existing

words are introduced in new context. For example, most newspapers and TV channels use 9/11 to refer to the terrorist act on the World Trade Centre in the USA in 2004. When we process written text or spoken utterances, we have access to underlying mental representation. The only way a machine can learn the meaning of a specific word in a message is by considering its context, unless some explicitly coded general world or domain knowledge is available. The context of a word is defined by co-occurring words. It includes everything that occurs before or after a word. The frequency of a word being used in a particular sense also affects its meaning. The English word ‘while’ was initially used to mean ‘a short interval of time’. But now it is more in use as a conjunction. None of the usages of ‘while’ discussed in this chapter correspond to this meaning.

Idioms, metaphor, and ellipses add more complexity to identify the meaning of the written text. As an example, consider the sentence:

The old man finally kicked the bucket.

The meaning of this sentence has nothing to do with the words ‘kick’ and ‘bucket’ appearing in it.

Quantifier-scoping is another problem. The scope of quantifiers (the, each, etc.) is often not clear and poses problem in automatic processing. The ambiguity of natural languages is another difficulty. These go unnoticed most of the times, yet are correctly interpreted. This is possible because we use explicit as well as implicit sources of knowledge. Communication via language involves two brains not just one—the brain of the speaker/writer and that of the hearer/reader. Anything that is assumed to be known to the receiver is not explicitly encoded. The receiver possesses the necessary knowledge and fills in the gaps while making an interpretation. As humans, we are aware of the context and current cultural knowledge, and also of the language and traditions, and utilize these to process the meaning. However, incorporating contextual and world knowledge poses the greatest difficulty in language computing. An example of cultural impact on language is the representation of different shades of white in the Eskimo world. It may be hard for a person living in plain to distinguish among various shades. Similarly, to an Indian, the word ‘Taj’ may mean a monument, a brand of tea, or a hotel, which may not be so for a non-Indian. Let us now take a look at the various sources of ambiguities in natural languages.

The first level of ambiguity arises at the word level. Without much effort, we can identify words that have multiple meanings associated with

them, e.g., bank, can, bat, and still. A word may be ambiguous in its part-of-speech or it may be ambiguous in its meaning. The word ‘can’ is ambiguous in its part-of-speech whereas the word ‘bat’ is ambiguous in its meaning. We hardly consider all possible meanings of a word to get the correct one. A program on the other hand, must be explicitly coded to resolve each meaning. Hence, we need to develop various models and algorithms to resolve them. Deciding whether ‘can’ is a noun or a verb is solved by ‘part-of-speech tagging’ whereas identifying whether a particular use of ‘bank’ corresponds to ‘financial institution’ sense or ‘river bank’ sense is solved by ‘word sense disambiguation’. Part-of-speech tagging and ‘word sense disambiguation’ algorithms are discussed in Chapters 3 and 5 respectively.

A sentence may be ambiguous even if the words are not, for example, the sentence: ‘*Stolen rifle found by tree*’. None of the words in this sentence is ambiguous but the sentence is. This is an example of structural ambiguity. Verb sub-categorization may help to resolve this type of ambiguity but not always. Probabilistic parsing, which is discussed in Chapter 4, is another solution. At a still higher level are pragmatic and discourse ambiguities. Ambiguities are discussed in Chapter 5.

A number of grammars have been proposed to describe the structure of sentences. However, there are an infinite number of ways to generate them, which makes writing grammar rules, and grammar itself, extremely complex. On top of it, we often make correct semantic interpretations of non-grammatical sentences. This fact makes it almost impossible for grammar to capture the structure of all and only meaningful text.

1.5 LANGUAGE AND GRAMMAR

Automatic processing of language requires the rules and exceptions of a language to be explained to the computer. Grammar defines language. It consists of a set of rules that allows us to parse and generate sentences in a language. Thus, it provides the means to specify natural language. These rules relate information to coding devices at the language level—not at the world-knowledge level (Bharati et al. 1995). However, since world knowledge affects both the coding (i.e., words) and the coding convention (structure), this blurs the boundary between syntax and semantics. Nevertheless such a separation is made because of the ease of processing and grammar writing.

The main hurdle in language specification comes from the constantly changing nature of natural languages and the presence of a large number

of hard-to-specify exceptions. Several efforts have been made to provide such specifications, which has led to the development of a number of grammars. Main among them are transformational grammar (Chomsky 1957), lexical functional grammar (Kaplan and Bresnan 1982), government and binding (Chomsky 1981), generalized phrase structure grammar, transformational grammar (Chomsky 1957), dependency grammar, Paninian grammar, and tree-adjoining grammar (Joshi 1985). Some of these grammars focus on derivation (e.g., phrase structure grammar) while others focus on relationships (e.g., dependency grammar, lexical functional grammar, Paninian grammar, and link grammar). We discuss some of these in Chapter 2. The greatest contribution to grammar comes from Noam Chomsky, who proposed a hierarchy of formal grammar based on level of complexity. These grammars use phrase structure rules (or rewrite rules). The term ‘generative grammar’ is often used to refer to the general framework introduced by Chomsky. Generative grammar basically refers to any grammar that uses a set of rules to specify or generate all and only grammatical (well-formed) sentences in a language. Chomsky argued that phrase structure grammars are not adequate to specify natural language. He proposed a complex system of transformational grammar in his book on *Syntactic Structures* (1957), in which he suggested that each sentence in a language has two levels of representation, namely, a deep structure and a surface structure (See Figure 1.1). The mapping from deep structure to surface structure is carried out by transformations. In the following paragraphs, we introduce transformational grammar.

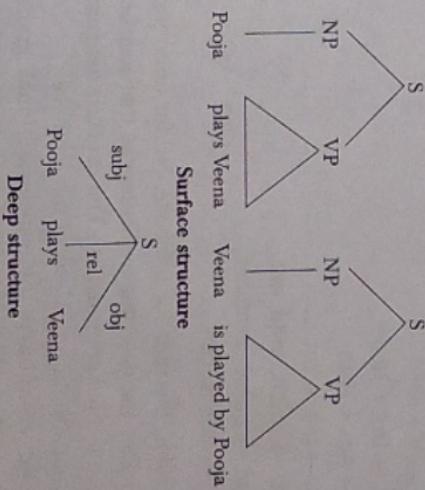


Figure 1.1 Surface and deep structures of sentence

Transformational grammar was introduced by Chomsky in 1957.

Chomsky argued that an utterance is the surface representation of a 'deeper structure' representing its meaning. The deep structure can be transformed in a number of ways to yield many different surface-level representations. Sentences with different surface-level representations having the same meaning, share a common deep-level representation.

Chomsky's theory was able to explain why sentences like

Pooja plays veena. (1.4a)

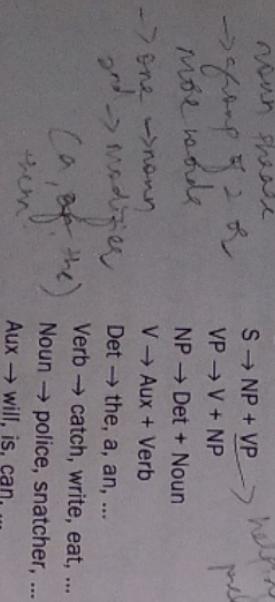
Veena is played by Pooja. (1.4b)

have the same meaning, despite having different surface structures (roles of subject and object are inverted). Both the sentences are being generated from the same 'deep structure' in which the deep subject is Pooja and the deep object is the veena.

Transformational grammar has three components:

1. Phrase structure grammar
2. Transformational rules
3. Morphophonemic rules—These rules match each sentence representation to a string of phonemes.

Each of these components consists of a set of rules. Phrase structure grammar consists of rules that generate natural language sentences and assign a structural description to them. As an example, consider the following set of rules:



In these rules, S stands for sentence, NP for noun phrase, VP for verb phrase, and Det for determiner. Sentences that can be generated using

these rules are termed grammatical. The structure assigned by the grammar is a constituent structure analysis of the sentence.

The second component of transformational grammar is a set of transformation rules, which transform one phrase-maker (underlying) into another phrase-marker (derived). These rules are applied on the terminal

string generated by phrase structure rules. Unlike phrase structure rules, transformational rules are heterogeneous and may have more than one symbol on their left hand side. These rules are used to transform one surface representation into another, e.g., an active sentence into passive one. The rule relating active and passive sentences (as given by Chomsky) is

$$NP_1 - Aux - V - NP_2 \rightarrow NP_2 - Aux + be + en - V - by + NP_1 \quad (1.5)$$

This rule says that an underlying input having the structure $NP - Aux - V - NP$ can be transformed to $NP - Aux + be + en - V - by + NP$. This transformation involves addition of strings 'be' and 'en' and certain rearrangements of the constituents of a sentence. Transformational rules can be obligatory or optional. An obligatory transformation is one that ensures agreement in number of subject and verb, etc., whereas an optional transformation is one that modifies the structure of a sentence while preserving its meaning. Morphophonemic rules match each sentence representation to a string of phonemes.

Consider the active sentence:

The police will catch the snatcher.

The application of phrase structure rules will assign the structure shown in Figure 1.2 to this sentence.

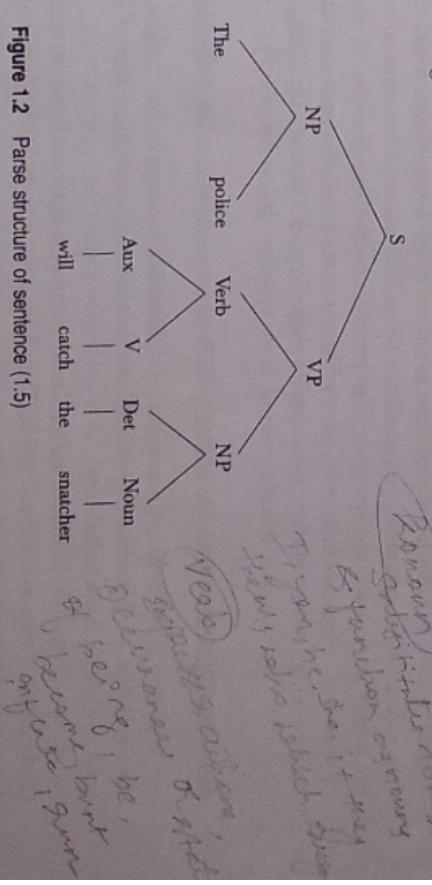


Figure 1.2 Parse structure of sentence (1.5)

The passive transformation rules will convert the sentence into:

The + culprit + will + be + en + catch + by + police (Figure 1.3).

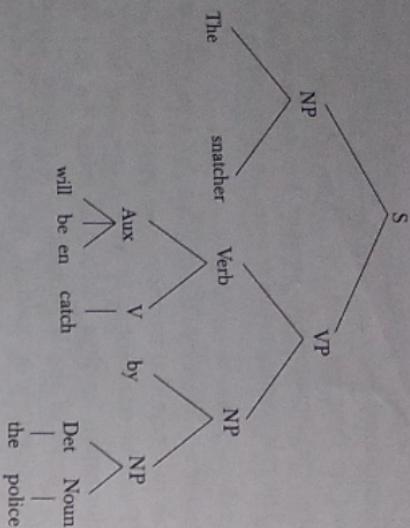


Figure 1.3 Structure of sentence (1.5) after applying passive transformations

Another transformational rule will then reorder ‘en + catch’ to ‘catch + en’ and subsequently one of the morphophonemic rules will convert ‘batch + en’ to ‘caught’. In general, the noun phrase is not always as simple as in sentence (1.5). It may contain other embedded structures, such as adjectives, modifiers, relative clause, etc. Long distance dependencies are other language phenomena that cannot be adequately handled by phrase structure rules. Long distance dependency refers to syntactic phenomena where a verb and its subject or object can be arbitrarily apart. The problem in the specification of appropriate phrase structure rules occurs because these phenomena cannot be localized at the surface structure level (Joshi and Vijayshankar 1989). Wh-movement¹ are a specific case of these types of dependencies.

1.6 PROCESSING INDIAN LANGUAGES

There are a number of differences between Indian languages and English. This introduces differences in their processing. Some of these differences are listed here.

- Unlike English, Indic scripts have a non-linear structure.
- Unlike English, Indian languages have SOV (Subject-Object-Verb) as the default sentence structure.

Paninian grammar provides a framework for Indian language models. These can be used for computation of Indian languages. The grammar focuses on extraction of Karaka relations from a sentence. We talk about the details of modelling in Chapter 2. A parsing framework based on Paninian grammar is introduced in Chapter 4 and issues involved in Indian language translation (using Paninian grammar theory) are discussed in Chapter 8.

1.7 NLP APPLICATIONS

Machine translation is the first application area of NLP. It involves the complete linguistic analysis of a natural language sentence, and linguistic generation of an output sentence. It is one of the most comprehensive and most challenging tasks in the area (AI-complete). However, the recent dramatic progress in the field of NLP has found interesting applications in information retrieval, information extraction, text summarization, etc. This book offers an extensive coverage of these recent applications, and also of traditional ones like machine translation and natural language generation. The focus has been on bridging the gap between theory and practice rather than on offering a gamut of linguistic, psychological, and computational theories.

¹It refers to a syntactic phenomenon in which interrogative words, called wh-words, appear at the beginning of the sentence. For example, when the direct object of the verb ‘read’ in the sentence ‘She is reading a book’ is replaced with a wh-word, the sentence becomes ‘What is she reading?’ instead of ‘She is reading what?’.

The applications utilizing NLP include the following.

Machine Translation

This refers to automatic translation of text from one human language to another. In order to carry out this translation, it is necessary to have an understanding of words and phrases, grammars of the two languages involved, semantics of the languages, and world knowledge.

Speech Recognition

This is the process of mapping acoustic speech signals to a set of words. The difficulties arise due to wide variations in the pronunciation of words homonym (e.g. dear and deer) and acoustic ambiguities (e.g., in the request and interest).

Speech Synthesis

Speech synthesis refers to automatic production of speech (utterance of natural language sentences). Such systems can read out your mails on telephone, or even read out a storybook for you. In order to generate utterances, text has to be processed. So, NLP remains an important component of any speech synthesis system.

Natural Language Interfaces to Databases

Natural language interfaces allow querying a structured database using natural language sentences.

Information Retrieval

This is concerned with identifying documents relevant to a user's query. NLP techniques have found useful applications in information retrieval. Indexing (stop word elimination, stemming, phrase extraction, etc.), word sense disambiguation, query modification, and knowledge bases have also been used in IR system to enhance performance, e.g., by providing methods for query expansion. WordNet, LDOCE (*Longman Dictionary of Contemporary English*) and Roget's *Thesaurus* are some of the useful lexical resources for IR research.

SysTran (System Translation)

The first SysTran machine translation system was developed in 1969 for Russian–English translation. SysTran also provided the first on-line machine translation service called Babel Fish, which is used by AltaVista search engines for handling translation requests from users.

TALM METEO

This is a natural language generation system used in Canada to generate weather reports. It accepts daily weather data and generates weather reports in English and French.

library—an information extraction system identifies a subset of information within a document that fits the pre-defined template.

Question Answering

Given a question and a set of documents, a question answering system attempts to find the precise answer, or at least the precise portion of text in which the answer appears. This is unlike an IR system, which returns whole document that seems relevant to the user's query. A question answering system is different from an information extraction system in that the content that is to be extracted is unknown. In general, a question answering system benefits from having an information extraction system to identify entities in the text. A question answering system requires more NLP than an information retrieval system or an information extraction system. It requires not only precise analysis of questions and portions of texts but also semantic as well as background knowledge to answer certain type of questions.

Text Summarization

This deals with the creation of summaries of documents and involves syntactic, semantic, and discourse level processing of text.

1.8 SOME SUCCESSFUL EARLY NLP SYSTEMS

ELIZA (Weizenbaum 1966)

ELIZA is one of the earliest natural language understanding programs. It uses syntactic patterns to mimic human conversation with the user. Here is a sample conversation.

Eliza: Hello. I am ELIZA. How may I help you?

User: I am feeling a little bit sleepy.

Eliza: How long have you been feeling a little bit sleepy?

User: For almost half an hour.

Eliza: Please go on.

***SHREDLU*(Winograd 1972)**

This is a natural language understanding system that simulates actions of a robot in a block world domain. It uses syntactic parsing and semantic reasoning to understand instructions. The user can ask the robot to manipulate the blocks, to tell the blocks configurations, and to explain its reasoning.

***LUNAR*(Woods 1977)**

This was an early question answering system that answered questions about moon rock.

1.9 INFORMATION RETRIEVAL

The availability of a large amount of text in electronic form has made it extremely difficult to get relevant information. Information retrieval systems aim at providing a solution to this.

The term ‘information’ should not be confused with the term ‘entropy’ (numerical measure of the uncertainty of an outcome) as it is used in communication theory. Information is being used here to reflect ‘subject matter’ or the ‘content’ of some text. We are not interested in ‘digital communication’, where bits and bytes are the information carriers. Instead our focus is on the communication taking place between human beings as expressed through natural languages. Information is always associated with some data (text, number, image, and so on); we are concerned with text only. Hence, we consider words as the carriers of information and written text as the message encoded in natural language.

As a cognitive activity, the word ‘retrieval’ refers to operation of accessing information from memory. We use the word ‘retrieval’ to refer to the operation of accessing information from some computer-based representation. Retrieval of information thus requires information to be processed and stored. Not all the information represented in computable form is retrieved. Instead, only the information relevant to the needs expressed in the form of query is located. In order to get this relevance, the stored and processed information needs to be compared against query representation. Information retrieval (IR) deals with all these facets. It is concerned with the organization, storage, retrieval, and evaluation of information relevant to the query.

Information retrieval deals with unstructured data. The retrieval is performed based on the content of the document rather than on its structure. The IR systems usually return a ranked list of documents. The IR components have been traditionally incorporated into different types

of information systems including database management systems, bibliographic text retrieval systems, question answering systems, and more recently in search engines.

Current approaches for accessing large text collections can be broadly classified into two categories. The first category consists of approaches that construct topic hierarchy, e.g., Yahoo. This helps the user locate documents of interest manually by traversing the hierarchy. However, it requires manual classification of new documents within the existing taxonomy. This makes it cost ineffective and inapplicable due to rapid growth of documents on the Web. The second category consists of approaches that rank the retrieved documents according to relevance. We discuss various IR models that support ranked retrieval in Chapter 9.

***Major Issues in Information Retrieval*(Siddiqui 2006)**

There are a number of issues involved in the design and evaluation of IR systems, which are briefly discussed in this section. The first important point is to choose a representation of the document. Most human knowledge is coded in natural language, which is difficult to use as knowledge representation language for computer systems. Most of the current retrieval models are based on keyword representation. This representation creates problems during retrieval due to polysemy, homonymy, and synonymy. Polysemy involves the phenomenon of a lexeme with multiple meaning. Homonymy is an ambiguity in which words that appear the same have unrelated meanings. Ambiguity makes it difficult for a computer to automatically determine the conceptual content of documents. Synonymy creates problem when a document is indexed with one term and the query contains a different term, and the two terms share a common meaning. Another problem associated with keyword-based retrieval is that it ignores semantic and contextual information in the retrieval process. This information is lost in the extraction of keywords from the text and cannot be recovered by the retrieval algorithms.

A related issue is that of inappropriate characterization of queries by the user. There can be many reasons for the vagueness and inaccuracy of the user’s queries, say for instance, her lack of knowledge of the subject or even the inherent vagueness of the natural language. The user may fail to include relevant terms in the query or may include irrelevant terms. Inappropriate or inaccurate queries lead to poor retrieval performance. The problem of ill-specified query can be dealt with by modifying or expanding queries. An effective technique based on user-interaction is relevance feedback which modifies queries based on the feedback provided by the user on initial retrieval.

In order to satisfy the user's request, an IR system matches document representation with query representation. Matching query representation with that of the document is another issue. A number of measures have been proposed to quantify the similarity between a query and the document to produce a ranked list of results. Selection of the appropriate similarity measure is a crucial issue in the design of IR systems.

Evaluating the performance of IR systems is also a major issue. There are many aspects of evaluation, the most important being the effectiveness of the system. Recall and precision are the most widely used measures of effectiveness.

As the major goal of IR is to search a document in a manner relevant to the query, understanding what constitutes relevance is also an important issue. Relevance is subjective in nature (Saracevic 1991). Only the user can tell the true relevance; it is not possible to measure this 'true relevance'. One may however, define the degree of relevance. Relevance has been considered as a binary concept, whereas it is in fact a continuous function (a document may be exactly what the user wants or it may be closely related). Current evaluation techniques do not support this continuity as it is quite difficult to put into practice. A number of relevance frameworks have been proposed (Saracevic 1996). These include the system, communication, psychological, and situational frameworks. The most inclusive is the situational framework, which is based on the cognitive view of the information seeking process and considers the importance of situation, context, multi-dimensionality, and time. A survey of relevance studies can be found in Mizzaro (1997). Most of the evaluations of IR systems have so far been done on document test collections with known relevance judgments.

The size of document collections and the varying needs of users also complicate text retrieval. Some users require answers of limited scope, while others require documents with a wider scope. These differing needs can require different and specialized retrieval methods. However, these are research issues and have not been dealt with in this book.

SUMMARY

- Language is the primary means of communication used by humans.
- Natural language processing is concerned with the development of computational models of aspects of human language processing.
- Theoretical linguists are mainly interested in providing a description of the structure and semantics of natural language, whereas

computational linguists deal with the study of language from a computational point of view.

- Historically, there have been two major approaches to natural language processing, namely rationalist approach and empiricist approach.
- The highly ambiguous and vague nature of natural language makes it difficult to create a representation amenable to computing.

REFERENCES

- Bach, Kent, 2002, *Meaning and Truth*, J. Kern Campbell, M. O'Rourke, and D. Shei (Eds.), Seven Bridges Press, New York, pp. 284–92.
- Chomsky, Noam, 1957, *Syntactic Structures*, Mouton, The Hague.
- , 1981, *Lectures on Government and Binding*, Foris Publications, Dordrecht, The Netherlands.
- Joshi, Aravind K., 1985, 'Tree adjoining grammar: How much sensitivity is required to provide reasonable structural description,' *Natural Language Parsing*, D. Dowty, L. Karttunen, and A. Zwicky (Eds.), Cambridge University Press, Cambridge.
- Joshi, Aravind K. and K. Vijayshankar, 1989, 'Treatment of long distance dependencies in LFG and TAG: functional uncertainty in LFG is a corollary in TAG,' *Proceedings of the 27th Annual Meeting on Association for Computational Linguistics*, Vancouver, British Columbia, pp. 220–27.
- Kaplan, R.M. and Joan Bresnan, 1982, 'Lexical functional grammar: A formal system for grammatical representation,' *The Mental Representation of Grammatical Relations*, Joan Bresnan (Ed.), MIT Press, Cambridge.
- Lenat, D.B., M. Prakash, and M. Shepherd, 1986, 'Cyc: using common sense knowledge to overcome brittleness and knowledge acquisition bottlenecks,' *AI Magazine*, 6(4).
- Mizzaro, S., 1997, 'Relevance: the whole history,' *Journal of the American Society for Information Science*, 48(9), pp. 810–32.
- Saracevic, T., 1991, 'Individual differences in organizing, searching and retrieving information,' *Proceedings of the 54th Annual Meeting of the American Society for Information Science (ASIS)*, pp. 82–86.
- , 1996, 'Relevance reconsidered,' *Proceedings of CoLIS 2, Second International Conference on Conceptions of Library and Information Science: Integration in Perspective*, P. Ingwersen and N.O. Pors (Eds.), The Royal School of Librarianship, Copenhagen, pp. 201–18.
- Siddiqui, Tanveer, 2006, 'Intelligent techniques for effective information retrieval: a conceptual graph-based approach,' *Ph.D. Thesis*, J.K. Institute of Applied Physics, Deptt. of Electronics and Communication, University of Allahabad.

Weizenbaum, R., 1966, 'ELIZA—A computer program for the study of natural language communication between man and machine,' *Communications of the ACM*, 9(1).

Winograd, Terry, 1972, *Understanding Natural Language*, Academic Press, New York.

Woods, William, 1977, 'Lunar Rocks in Natural English: Explorations in Natural Language Question-answering,' *Linguistic Structures Processing*, A. Zampoli [Ed.], Elsevier, North Holland.

EXERCISES

1. Differentiate between the rationalist and empiricist approaches to natural language processing.
2. List the motivation behind the development of computational models of languages.
3. Briefly discuss the meaning components of a language.
4. What makes natural language processing difficult?
5. What is the role of transformational rules in transformational grammar? Explain with the help of examples.

CHAPTER 2

LANGUAGE MODELLING

CHAPTER OVERVIEW

The domain of language is quite vast. It presents an almost infinite number of sentences to the reader (or computer). To handle such a large number of sentences, we have to create a model of the domain, which can subsequently be simplified and handled computationally. A number of language models have been proposed. We introduce some of these models in this chapter. To create a general model of any language is a difficult task. There are two approaches for language modelling. One is to define a grammar that can handle the language. The other is to capture the patterns in a grammar language statistically. This chapter has a mixed approach. It gives a glimpse of both grammar-based model and statistical language model. These include lexical functional grammar, government and binding, Pāṇini grammar, and n-gram based model.

2.1 INTRODUCTION

Why and how do we model a language? This question has been discussed by linguists since 500 BC. Computational linguists also have to confront this question. It is obvious that our purpose is to understand and generate natural languages from a computational viewpoint. One approach can be to just take a language, try to understand every word and sentence of it, and then come to a conclusion. This approach has not succeeded as there are difficulties at each stage, which we will understand as we go through this book. An alternative approach is to study the grammar of various languages, compare them, and if possible, arrive at reasonable models that facilitate our understanding of the problem and designing of natural-language tools.

A model is a description of some complex entity or process. A language model is thus a description of language. Indeed, natural language is a complex entity and in order to process it through a computer-based program, we need to build a representation (model) of it. This is known

This section focuses on lexical functional grammar (LFG), generative lexico-grammatical approaches to create statistical models of language and grammar language in a grammatical and rule-based format. It also introduces the unification grammar (UG) and introduces various approaches to understand these structures on lexical functional grammar (LFG), generative lexico-grammatical approaches to create statistical models of language and grammar.

2.2.1 Generative Grammars

Chomsky (1956) described classes of grammars in a hierarchical manner where the top layer contained the grammars represented by this sub-class. Hence, Type 0 (unrestricted) grammars are represented by this sub-class. Semisyntax grammar, which in turn contains Type 0 (context-free grammars) and that again contains Type 3 grammars (regular grammars). Although this relationship has been given for classes of formal grammars, it can be extended to descriptive grammars at various levels, such as in a class (class membership) relationship.

2.2.2 Hierarchical Grammar

Linguage is a relation between the sound (or the written text) and its meaning. Thus, any model of a language should also deal with the meaning of the utterances. As seen earlier, we can have a perfect grammar which matches all the sentences. In this chapter, we will assume that grammars are a type of language that can generate sentences.

In 1957, in his book on *Syntactic Structures*, Noam Chomsky wrote that we can generate sentences in a language if we know a collection of words and rules in that language. Only those sentences that can be generated and rules in the language. Of course, we are limiting only the same idea can be used to model a language. If we have a complete set of rules that can generate sentences in a language, those sentences provide a model of that language. Of course, we are limiting only the same idea can be used to model a language. If we have a complete set of rules that can generate sentences in a language, those sentences provide a model of that language.

2.2.1 Generative Grammars

dominant approaches to create stable

This section focuses on lexical microfunctional grammar (LMG), generated by three semantic grammars (G_S): government and binding (GB), and domain grammars (DG). LMG also introduces various approaches to microfunctions in a grammatical and rule-based format; it also introduces

Various computational grammars have been proposed and studied, e.g., transformational grammar (Chomsky 1957), lexical functional grammar (Kaplan and Bresnan 1982), government-and-binding (Chomsky 1981), generalized phrase structure grammar (Gazdar et al. 1985), dependency grammar, patterned grammar, and tree-adjoining grammar (Joshi 1985).

GRAMMAR-BASED LANGUAGE MODELS

Semantic language modeling is one of the fundamental tasks in many NLP applications, including speech recognition, spelling correction, handwriting recognition, and machine translation. It has now found applications in information retrieval, text summarization, and question answering [3]. A number of statistical language models have been proposed in literature. The most popular of these are the n-gram models [4]. We discuss this model in Section 2.3. The following section discusses various grammar-based models.

The statistical approach creates a language model by training it from a corpus, in order to capture regularities of a language. The training corpus needs to be sufficiently large. Rosenfeld (1994) pointed out: "The statistical approach creates a language model by training it from a corpus, in order to capture regularities of a language. The training corpus needs to be sufficiently large. Rosenfeld (1994) pointed out: "Statistical language modeling (SLM) is the attempt to capture performance of various natural language applications.

A grammar-based approach uses the grammar of a language to create its model. It attempts to represent the syntactic structure of language model. Grammar consists of hand-coded rules defining the structure and ordering of various constituents appearing in a linguistic unit (phrase, sentence, etc.). For example, sentence usually consists of noun phrase and a verb phrase. The grammar-based approach attempts to utilize this structure and also the relationships between these structures.

Language modeling can be viewed either as a problem of grammar inference or a problem of probability estimation. A grammar-based language modeling approach tries to distinguish a grammatical language model from a non-grammatical one by means of statistical measures, usually a maximum likelihood estimation. These two viewpoints disagree: models that empirically fit a set of training data do not necessarily fit a different set of test data. This is because we do not have a probability measure, but only a language organization of language models.

The X theory (pronounced X-bar theory) is one of the central concepts in GB. Instead of defining everel phrasal structures and the sentence structure with separate sets of rules, X theory defines them both as maximal projections of some head. In this manner, the entities as phrases of noun (N), verb (V), adjective (A), and preposition phrase (PP) are maximal projections of the head phrase (AP), and prepositional phrase (NP), verb phrase (VP), adjective phrase (AP), and complement phrase (CP), even the sentence phrases (where $X = (N, V, A, P)$). Not only that, even the corresponding respects, and can be represented as head X of their corresponding (P).

Figure 2.6 Organization of GB (adapted from Peter Selis 1985)

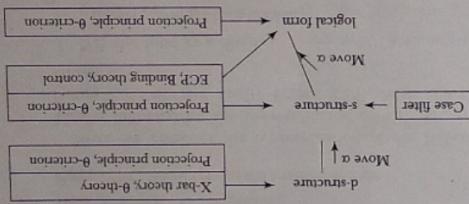


Figure 2.5 showed the application of various theories and principles at three different levels of representation in GB. Figure 2.6 mentions these explicitly to understand the organization of GB.

In LFI, the interpretation is that most travellers visit the same two countries (say, India and China). In L2F, when we move [most travellers] outside the scope of the sentence, the interpretation can be that most travellers visit two countries, which may be different travellers. One of the grammar which prohibits certain combinations and movements; it is the imperative Move a can move anything to any possible position. Thus, GB basically the formulation of theories or principles which create constraints to disallow the construction of ill-formed sentences. To account for cross-lingual constraints of similar type, GB can specify that a constituent cannot be moved from position X , where X can have value X_1 in one arrangement, X_2 in another, and so on. These rules are so general and amaguage-independent that language-particular details of description probably go uncharted in GB. (Sells 1985).

Language Modelling 27

Example 2.2 Consider the sentence:

Two countries are visited by [NP most travellers].

This two possible logical forms are:

LF1: [NP Two countries are visited by [NP most travellers]]

LF2: Applying Move α

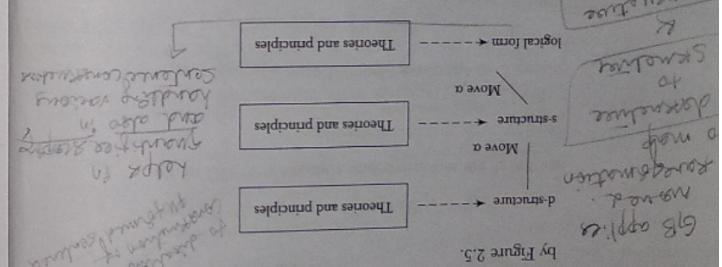
It is clear that LF2 is more appropriate because it reflects the intended meaning of the sentence.

The GB considers all three levels of representations (d-, s-, and LF) as syntactic, and GB is also related to semantics or semantic interpretation mechanisms. However, GB applies the same Move α transformation to map d-levels to s-levels or s-levels to LF level. LF level helps us in understanding or integrating various semantic constructions such as passive scopings and also in handling core constructions. An example of LF representation may be helpful.

C-command and government, case theory, unifying category principle (TCP), C-complexes (Selkirk 1985) applied at various levels of its representations and the transformation rule, Move α. Before elaborating on these modules— which include X-bar theory, projection principle, θ-theory and φ-theory— we discuss the general characteristics of GB.

The GB consists of a set of modules that contain constituents and binding theory—we discuss the general characteristics of GB.

89



Components of GB

Government and binding (GB) comprises a set of theories which map the structures of the phonetic form (LF). A generative form (GF) structure and its syntactic structure are mapped into the GB structure. Hence, in its simplest form GB can be represented as follows:

26 Natural Language Processing and Information Retrieval

5. S: that she ate the food in a dhaba
 [S [Ccomp that] [S [D_d, she] [INFL past]]]_p ate the food in a dhaba]]

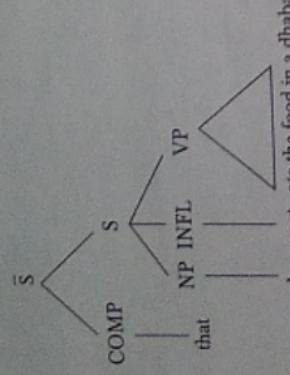


Figure 2.7(f) Maximal projection of sentence structure

As shown in Figure 2.7(f), the sentence is considered to be the head of INFL and the projection of sentence is denoted by \bar{S} , which has the specifier as complementizer (COMP).

Sub-categorization It is to be noted that GB does not consider traditional phrase structure as an appropriate device for defining language constructs. It places the burden of ascertaining well-formedness to sub-categorization frames of heads. In principle, any maximal projection can be the argument of a head, but sub-categorization is used as a filter to permit various heads to select a certain subset of the range of maximal projections. For example, we know that the verb 'eat' can sub-categorize for NP, whereas the verb 'sleep' cannot. Hence, 'ate food' is well-formed, but the sentence "slept the bed" is not. GB claims that defining phrase structures as head projections and sub-categorization helps ensure well-formed structures, even at the sentence level.

As 'verb' is not the head of the sentence, it cannot sub-categorize for 'subject NP', while it can perfectly sub-categorize for 'object NP'. This explains, to certain extent, the 'subject/object' asymmetry (Sells 1985).

Projection Principle

The projection principle, a basic notion in GB, places a constraint on the three syntactic representations and their mapping from one to the other. The principle states that representations at all syntactic levels (i.e., d-level, s-level, and LF level) are projections from the lexicon. Thus, lexical

*will precede the incorrect sentences.

properties of categorical structure (sub-categorization) must be observed at each level.

This can be understood with an example. Suppose 'the object' is not present at d-level, then another NP cannot take this position at s-level. This principle, in conjunction with the possibility of presence of empty category and other theories (like Binding Theory) ensures correct movement and well-formed structure.

Theta Theory (θ -theory) or The Theory of Thematic Relations

As discussed earlier, 'sub-categorization' only places a restriction on syntactic categories which a head can accept. GB puts another restriction on the lexical heads through which it assigns certain roles to its arguments. These roles are pre-assigned and cannot be violated at any syntactical level as per the projection principle. These role assignments are called theta-roles and are related to 'semantic-selection'.

Theta-role and Theta-criterion There are certain thematic roles from which a head can select. These are called θ -roles and they are mentioned in the lexicon, say for example the verb 'eat' can take arguments with θ -roles '(Agent, Theme)'. Agent is a special type of role which can be assigned by a head to outside arguments (external arguments) whereas other roles are assigned within its domain (internal arguments). Hence in 'Mukesh ate food', the verb 'eat' assigns the 'Agent' role to 'Mukesh' (outside VP) and 'Theme' (or 'patient') role to 'food'. As roles are assigned based on the syntactic positions of the arguments, it is important that there should be a match between the number of roles and number of arguments as depicted by θ -criterion.

Theta-Criterion states that 'each argument bears one and only one θ -role, and each θ -role is assigned to one and only one argument' (Sells 1985). Thus, each argument will have a unique θ -role and cannot be moved to a position where it may acquire another θ -role.

In GB, d-structure is conceived as some kind of 'pure' representation of arguments and hence, θ -roles are assigned at d-level only, whereas theta-criterion is applied at all the three levels, as shown in Figure 2.6.

C-command and Government

As 'Government' is a special case of 'C-command', we will first define C-command.

C-command C-command defines the scope of maximal projection. It is a basic mechanism through which many constraints are defined on Move

α . If any word or phrase (say α or β) falls within the scope of α and is determined by a maximal projection, we say that it is dominated by the maximal projection. Now, if there are two structures α and β related in such a way that 'every maximal projection dominating α dominates β ', we say that α C-commands β , and this is the necessary and sufficient condition (iff) for C-command.

The definition of C-command does not include all maximal projections dominating β , only those dominating α . If we put this extra constraint, it becomes a kind of mutual C-command (Sells 1985), called government.

Government

- α governs β iff:
- α C-commands β

α is an X (head, e.g., noun, verb, preposition, adjective, and inflection), and every maximal projection dominating β dominates α .

Thus no maximal projection can intervene between the governor and governed. In GB literature, this has been stated as: 'Maximal projections are barriers to government.'

Movement, Empty Category, and Co-indexing

Briefly let us discuss Move α . In GB, Move α is described as 'move anything anywhere', though it provides restrictions for valid movements.

In GB, the active to passive transformation is the result of NP movement as shown in sentence (2.3). Another well-known movement is the wh-movement, where wh-phrase is moved as follows.

- What did Mukesh eat? (2.3)
[Mukesh INFL eat what]

As discussed in the projection principle, lexical categories must exist at all the three levels. This principle, when applied to some cases of movement leads to the existence of an abstract entity called empty category. In GB, there are four types of empty categories, two being empty NP positions called wh-trace and NP trace, and the remaining two being pronouns called small 'pro' and big 'PRO'. This division is based on two properties—anaphoric (+a or -a) and pronominal (+p or -p).

- Wh-trace -a, -p
- NP-trace +a, -p
- small 'pro' -a, +p
- big 'PRO' +a, +p

Co-indexing is the indexing of the subject NP and AGR (agreement) at d-structure which are preserved by Move α operations at s-structure.

When an NP-movement takes place, a trace of the movement is created by having an indexed empty category (e_i) from the position at which the movement began to the corresponding indexed NP, i.e. NP_i . All A-positions (argument positions) at s-level are also freely indexed. These categories and indices are used to define Binding Theory.

It is interesting to note that for defining constraints to movement, the theory identifies two positions in a sentence. Positions assigned θ -roles are called θ -positions, while others are called $\bar{\theta}$ -positions.

In a similar way, core grammatical positions (where subject, object, indirect object, etc., are positioned) are called A-positions (arguments positions), and the rest are called \bar{A} -positions.

Binding Theory

Binding is defined by Sells (1985) as follows:

- α binds β iff
- α C-commands β , and
- α and β are co-indexed

As we noticed in sentence (2.1),

- [e_i INFL kill Mukesh]
- [Mukesh, was killed (by e_j)]
- Mukesh was killed.

Empty clause (e_i) and Mukesh (NP_i) are bound. This theory gives a relationship between NPs (including pronouns and reflexive pronouns).

Now, binding theory can be given as follows:

- (a) An anaphor (+a) is bound in its governing category.
- (b) A pronominal (+p) is free in its governing category.
- (c) An R-expression (-a, -p) is free.

This theory applies to binding at A-positions. Governing category is the local domain (the smallest only) NP or S containing it (G or p or R-expression) and its governor.

Example 2.4

- A: Mukesh_i knows himself,
- B: Mukesh_i believes that Amrita knows him,
- C: Mukesh believes that Amrita_j knows Nupur_k

Similar rules apply on empty categories also:

- NP-trace: +a, -p: Mukesh_i was killed e_i
- wh-trace: -a, -p: Who_i does he_j like e_i

Empty Category Principle (ECP)

We have already defined 'government'. Now, let us define 'proper government':

- α properly governs β iff:
- α governs β and α is lexical (i.e. N, V, A, or P) or
- α locally A-binds β

The ECP says 'A trace must be properly governed'.

This principle justifies the creation of empty categories during NP-trace and wh-trace and also explains the subject/object asymmetries to some extent. As in the following sentences:

- What, do you think that Mukesh ate e_i ?
- What, do you think Mukesh ate e_i ?

Bounding and Control Theory

There are many other types of constraints on Move α . It is not possible to explain all of them here, for details, see Peter Sells (1985).

In English, the long distance movement for complement clause can be explained by bounding theory if NP and S are taken to be bounding nodes. The theory says that the application of Move α may not cross more than one bounding node. The theory of control involves syntax, semantics, and pragmatics. As stated previously, the empty category 'PRO' (+a, +p) behaves as an anaphor sometimes, when it is the subject of the clausal complement to verbs such as decide and try. However, it behaves as pronoun with some other verbs.

Case Theory and Case Filter

In GB, case theory deals with the distribution of NPs and mentions that each NP (with the possible exception of a few empty categories) must be assigned a case. In English, we have the nominative, objective, genitive, etc., cases, which are assigned to NPs at particular positions. Indian languages are rich in case-markers, which are carried even during movements.

Case Filter An NP is ungrammatical if it has phonetic content or if it is an argument (with the exception of big 'PRO') and is not case-marked. Phonetic content here, refers to some physical realization, as opposed to empty categories. Thus, case filters restrict the movement of NP at a position which has no case assignment. It works in a manner similar to that of the θ -criterion.

In short, GB presents a model of the language which has three levels of syntactic representation. It assumes phrase structures to be the maximal

projection of some lexical head and in a similar fashion, explains the structure of a sentence or a clause. It assigns various types of roles to these structures and allows them a broad kind of movement called Move α . It then defines various types of constraints which restrict certain movements and justifies others. GB gives a new insight for the modelling of languages, although the Chomskian Minimalist Programme has superseded GB.

2.2.4 Lexical Functional Grammar (LFG) Model

This section presents those features of LFG that throw a light on language modelling. For the details of lexical functional grammar, readers are encouraged to seek Darlymple et al. (1995).

Unlike GB, LFG represents sentences at two syntactic levels—constituent structure (c-structure) and functional structure (f-structure). Based on Woods' *Augmented Transition Networks* 1970, which used phrase structure trees to represent the surface structure of sentences and the underlying predicate-argument structure; Kaplan (1975a, b) proposed a concrete form for the register names and values (used in ATN implementation), which became the functional structures in LFG. On the other hand, Bresnan (1976a, 1977) was more concerned with the problem of explaining some linguistic issues, such as active/passive and dative alternations, in transformational approach. She proposed that such issues can be dealt with by using lexical redundancy rules. The unification of these two diverse approaches (with a common concern) led to the development of the LFG theory, which was presented as *Lexical Functional Grammar: A Formal System for Grammatical Representation* in 1982.

The LFG is a formalism that is both computationally and linguistically motivated and provides precise algorithms for linguistic issues it can handle. The term 'lexical functional' is composed of two terms: the 'functional' part is derived from 'grammatical functions', such as subject and object, or roles played by various arguments in a sentence. The 'lexical' part is derived from the fact that the lexical rules can be formulated to help define the given structure of a sentence and some of the long distance dependencies, which is difficult in transformational grammars.

C-structure and f-structure in LFG

As LFG is aimed at providing exact computational algorithms, it provides well-defined objects called constituent structure (c-structure) and functional structure (f-structure). The c-structure is derived from the usual phrase and sentence structure syntax, as in CFG (discussed in Chapter 4). However,

Finally, the f-structure is the set of attribute-value pairs, represented as

subj	Pers	3	[Pers Num SG Gen FEM Case NOM Pred 'PRO']
	Num	SG	
	Gen	FEM	
obj	Case	NOM	[Pers Num PL Pred 'Star']
	Pers	3	
	Num	PL	
Pred	'see'	<(↑ subj) (↑ obj)>	

It is interesting to note that the final f-structure is obtained through the unification of various f-structures for subject, object, verb, complement etc. This unification is based on the functional specifications of the verb which predicts the overall sentence structure.

LFG requires that all possible structures corresponding to passive constructs, dative constructs, etc., must be specified. If the given sentence does not match the specifications, it is said to be ill-formed. LFG imposes three conditions on f-structure (Sells 1985).

Consistency In a given f-structure, a particular attribute may have a single value. Hence, while unifying two f-structures, if the attribute Num has value SG in one and PL in the other, it will be rejected.

Completeness A function is called governable if it appears within the Pred value of some lexical form, e.g., Subj, Obj, and Obj 2. Adjunct is not a governable function.

When an f-structure and all its subsidiary f-structures (as the value of any attribute of f-structure can again contain other f-structures) contain all the functions that their predicates govern, then and only then is the f-structure complete. For example, since the predicate 'see <(↑ Subj (↑ Obj)>' contains an object as its governable function, a sentence like 'He saw' will be incomplete.

Coherence Coherence maps the completeness property in the reverse direction. It requires that all governable functions of an f-structure, and all its subsidiary f-structures, must be governed by their respective predicates. Hence, in the f-structure of a sentence, an object cannot be taken if its verb does not allow that object. Thus, it will reject the sentence, 'I laughed a book.'

The completeness and coherence conditions are counterparts of θ-criterion in GB theory.

Lexical Rules in LFG

Different theories have different kinds of lexical rules and constraints for handling various sentence-constructs (active, passive, dative, causative, etc.). In GB, to express a sentence in its passive form, the verb is changed to its participial form and the ability of the verb to assign case and external (Agent) θ-role is taken away. In LFG, the verb is converted to the participial form, but the sub-categorization is changed directly. Consider the following example:

Active: Tara ate the food.

Passive: The food was eaten by Tara.

Active: ↑ Pred = 'eat <(↑ Subj) (↑ Obj)>'

Passive: ↑ Pred = 'eat <(↑ Obj) (↑ Subj)>'

Here, Obj_{ag} represents oblique agent phrase. Similar rules can be applied in active and dative constructs for the verbs that accept two objects.

Active: Tara gave a pen to Monika.

Passive: Tara gave Monika a pen.

Active: ↑ Pred = 'give <(↑ Subj) (↑ Obj) (↑ Obj)>'

Passive: ↑ Pred = 'give <(↑ Subj) (↑ Obj) (↑ Obj) (↑ Obj)_{go}>'

Here, Obj_{go} stands for oblique goal phrase. Similar rules are also applicable to the process of causativization. This can be seen in Hindi, where the verb form is changed as follows:

तारा लहसी
Taraa lahsii
Laugh
→ causativation
तारा लहसी
Taraa hansi
Tara laughed

Example 2.6

Active:

तारा हंसी
Taraa hansi
Tara laughed

Causative:
मोनिका ने तारा को हँसाया
Monika ne Tara ko hansayaa
Monika Subj Tara Obj laugh-cause-past
Monika made Tara to laugh.

Active: ↑ Pred = 'Laugh <(↑ Subj)>'

Causative: ↑ Pred = 'cause <(↑ Subj) (↑ Obj) (Comp)>'

Here, a new predicate is formed which causes the action and requires a new subject, while the old subject becomes the object of the new predicate (complement to infinitive and the old verb becomes VPs).

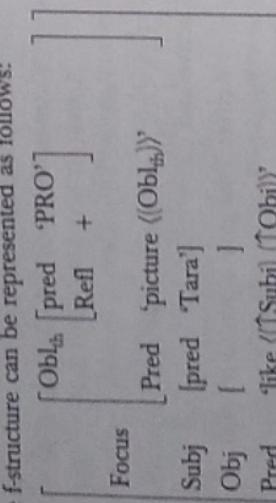
Long Distance Dependencies and Coordination

In GB, when a category moved, it creates an empty category. In LFG, unbounded movement and coordination is handled by the functional identity and by correlation with the corresponding f-structure. An example will better explain these ideas.

Example 2.7 Consider the wh-movement in the following sentence.

Which picture does Tara like—most?

The f-structure can be represented as follows:



The mechanism of handling these movements and coordination are not detailed here. The only aim is to highlight efforts and issues involved in modelling language.

2.2.5 Paninian Framework

Another very important model which has drawn much attention is the Paninian Grammar-based model (Kiparsky 1982, Bharti et al. 1995). Although *Paninian grammar* (PG) was written by Panini in 500 BC in Sanskrit (the original text being titled *Asthadhyayi*), the framework can be used for other Indian languages and possibly some Asian languages as well.

Unlike English, Asian languages are SOV (Subject-Object-Verb) ordered and inflectionally rich. The inflections provide important syntactic and semantic cues for language analysis and understanding. The Paninian framework takes advantage of these features. However, it should be noted that the research on this framework is still in progress and there are many complexities of Indian languages which are yet to be explained through this or other models. In this section, we briefly discuss some unique features of PG, to provide a glimpse of another potential model.

Some Important Features of Indian Languages

Indian languages have traditionally used oral communication for knowledge propagation. The purpose of these languages is to communicate ideas from the speaker's mind to the listener's mind. Such oral traditions have given rise to a morphologically rich language. Also, they are relatively word-order free. Some languages, like Sanskrit, have the flexibility to allow word groups representing subject, object, and verb to occur in any order. In others, like Hindi, we can change the position of subject and object. For example:

(a) मौं बच्चे को लाला देती है।

Maan Bachche ko khanaa detii hai
Mother child to food give—(s)

(b) बच्चे को मौं लाला देती है।

Bachche ko Maan khanaa detii hai
Child to mother food give—(s)
Mother gives food to the child.

The auxiliary verbs follow the main verb. In Hindi, they remain as separate words, whereas in south Indian (Dravidian) languages, they combine with the main verb. For example:

करता रहा है।
kartaah rahaah hai
doing been has
has been doing
eating

In Hindi, some verbs (main), e.g., give (देना), take (लेना), also combine with other verbs (main) to change the aspect and modality of the verbs.

Example 2.8

उसने आला आया।

Usne khaanaa khaayaa
He (Subj) food ate
He ate food
वह चल दिया।

He move given
He moved (started the action)

In Indian languages, the nouns are followed by postpositions instead of prepositions. They generally remain as separate words in Hindi, except in the case of pronouns, for example

रेखा के पिता उसके पिता

Rekha ke pita Uske pita

Father of Rekha Father of Rekha

In view of such differences between English (and English-like languages) and Indian languages, it is imperative that we find a new framework for handling Indian languages. Even among Indian languages, all features are not the same. As noted earlier, verb groups are formed differently in Indo-Aryan and Dravidian languages. Sanskrit is very different from the other Indian languages as it has five tenses and three numbers, and only one time aspect in each tense. Hence, the translation of 'He goes' and 'He is going' is the same in Sanskrit. Hindi is unique in the sense that it has no neuter gender. All nouns are categorized as feminine or masculine, and the verb form must have a gender agreement with the subject (sometimes with the object). Thus, we have

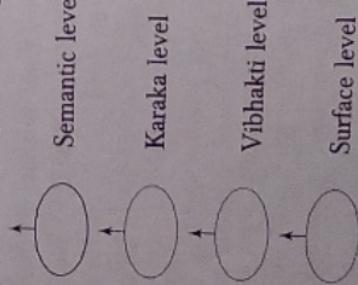
चाही खो जवी

Chaahee kho gayee
key lose (past)

The key was lost.

Layered Representation in PG

The GB theory represents three syntactic levels: deep structure, surface structure, and logical form (LF), where the LF is nearer to semantics. This theory tries to resolve all language issues at syntactic levels only. Unlike GB, Paninian grammar framework is said to be syntactico-semantic, that is, one can go from surface layer to deep semantics by passing through intermediate layers. Although all these layers are not named, as per Bharti, et al. (1995), the language can be represented as follows:



The surface and the semantic levels are obvious. The other two levels should not be confused with the levels of GB. Vibhakti literally means inflection, but here, it refers to word (noun, verb, or other) groups based either on case endings, or post-positions, or compound verbs, or main and auxiliary verbs, etc. Instead of talking about NP, VP, AP, PP, etc., word groups are formed based on various kinds of markers (including the absence of it or θ). These markers are language-specific, but all Indian languages (and possibly Asian languages as well) can be represented at the Vibhakti level.

Karaka (pronounced *Karakā*) literally means Case, and in GB, we have already discussed case theory, θ-theory, and sub-categorization, etc. Paninian Grammar has its own way of defining Karaka relations, which we discuss in the next section. These relations are based on the way the word groups participate in the activity denoted by the verb group. In this sense, it is semantic as well as syntactic. However, things are not so straightforward. Complexities arise because of the absence of inflections, multiple categories of the words, multiple meanings, and above all, the presence of a large number of exceptions. These exceptions are not only applicable on stated rules but also on future rules. Such forward and backward chaining makes actual implementation difficult.

As the purpose of these languages is to communicate, generally between one human and another, the resolution of ambiguities is a contentious issue, often left to the listener. Hence, there may not be any particular number of semantic levels. Multiple-meaning texts are abundant in Indian literature as seen in the hundreds of interpretations of the epics.

Karaka Theory

Karaka theory is the central theme of PG framework. Karaka relations are assigned based on the roles played by various participants in the main activity. These roles are reflected in the case markers and post-position markers (parsargs). These relations are similar to case relations in English, but the types of relations are defined in a different manner and the richness of the case endings found in Indian languages has been used to its advantage.

We will discuss the various Karakas, such as Karta (subject), Karma (object), Karana (instrument), Sampradana (beneficiary), Apadan (separation), and Adhikaran (locus). These descriptions are just examples and not a complete discussion of PG or Karaka theory. For details of Karaka theory, see Shastri (1973) and Vasu (1977).

To explain various Karaka relations, let us consider an example.

Consider the following Hindi sentence:

माँ बच्ची को औलन में हाथ से रोटी लिखाती है।

Maan bachi ko aangan mein haath se rotii khilaati hei

Mother child-to courtyard-in hand-by bread feed (s).

The mother feeds bread to the child by hand in the courtyard.

The first important Karak is subject, called 'Karta' in PG. Karta is defined as the noun group which is most independent (*svatantra* in Hindi). Karta has generally 'ne' or 'ø' case marker. It is an independent entity in the activity denoted by the main verb. As seen in GB also, verbs do not sub-categorize for subjects, although they assign a θ role to it. In sentence 2.5, 'maan' (mother) is the Karta. The concept of Karta is different from the 'agent' concept in the sense that Karta can also take up the role of experiencer, e.g.,

मुझसे रहा न जाया।

Mujhe raha na gayaa

Me hold -not -passive

I could not hold myself.

'Karma' is similar to object and is the locus of the result of the activity. In sentence (2.5), *rotii* (bread) is the Karma. As explained earlier, when the Karta is the experiencer, it (she) is also the locus of the result. Thus, the locus of the result is only when it is different from Karta termed Karma. Karma generally has 'ø' or 'KO' case marker. Another Karaka relation is 'Karan' (instrument), which is a noun group through which the goal is achieved. In sentence (2.5), *haath* (hand) is the Karan. It has the marker *dvara* (by) or *sz*. 'Sampradan' is the beneficiary of the activity, e.g., *bachchi* (child) in sentence (2.5). It takes the marker *ko* (to) or *ke lya* (for). 'Apadaan' denotes separation and the marker is attached to the part that serves as a reference point (being stationary), for example

माँ ने थाली से आला उपादान बच्चे को दिया।

Maan ne thaali se khana uthakar bachche ko diyaa

Mother-Karta plate from Apaadaan food taking-up child-to gave.

The mother gave food to the child taking it up from the plate.

Here *thaali* is the Apaadaan. 'Adhikaran' is the locus (support in space or time) of Karta or Karma. In sentence (2.5), *aangan* (courtyard) is the Adhikaran. As these six relations are not sufficient to capture all possible relations, various others such as 'Sambandh' (relation) and 'Tadarthyā' (purpose) have also been tried.

Issues in Paninian Grammar

The two problems challenging linguists are:

- (i) Computational implementation of PG, and
- (ii) Adaptation of PG-to Indian, and other similar languages.

An approach to implementing PG has been discussed in Bharati, et al. (1995). This is a multilayered implementation. The approach is named 'Utsarga-Apvida' (default-exception), where rules are arranged in multiple layers in such a way that each layer consists of rules which are in exception to rules in the higher layer. Thus, as we go down the layer, more particular information is derived. Rules may be represented in the form of charts (such as Karaka chart and Lakshan chart).

However, many issues remain unresolved, especially in cases of shared Karak relations. Another difficulty arises when mapping between the Vibhakti (case markers and post-positions) and the semantic relation (with respect to verb) is not one to one. Two different Vibhakti can represent the same relation, or the same Vibhakti can represent different relations in different contexts. The strategy to disambiguate the various senses of words, or word groupings, are still the challenging issues.

As the system of rules is different in different languages, the framework requires adaptations to tackle various applications in various languages. Only some general features of PG framework has been described here.

2.3 STATISTICAL LANGUAGE MODEL

A statistical language model is a probability distribution $P(s)$ over all possible word sequences (or any other linguistic unit like words, sentences, paragraphs, documents, or spoken utterances). A number of statistical language models have been proposed in literature. The dominant approach in statistical language modelling is the n -gram model.

2.3.1 n -gram Model

As discussed earlier, the goal of a statistical language model is to estimate the probability (likelihood) of a sentence. This is achieved by decomposing sentence probability into a product of conditional probabilities using the chain rule as follows:

$$\begin{aligned} P(s) &= P(w_1, w_2, w_3, \dots, w_n) \\ &= P(w_1) P(w_2|w_1) P(w_3|w_1 w_2) P(w_4|w_1 w_2 w_3) \dots \\ &\quad = \prod_{i=1}^n P(w_i|h_i) \end{aligned}$$

where h_i is history of word w_i defined as

$$w_1 w_2 \dots w_{i-1}$$

So, in order to calculate sentence probability, we need to calculate the probability of a word, given the sequence of words preceding it. This is not a simple task. An n -gram model simplifies the task by approximating the probability of a word given all the previous words by the conditional probability given previous $n-1$ words only.

$$P(w_i/h_i) \approx P(w_i/w_{i-n+1} \dots w_{i-1})$$

Thus, an n -gram model calculates $P(w_i/h_i)$ by modelling language as Markov model of order $n-1$, i.e., by looking at previous $n-1$ words only. A model that limits the history to the previous one word only, is termed a bi-gram ($n=1$) model. Likewise, a model that conditions the probability of a word to the previous two words, is called a tri-gram ($n=2$) model. Using bi-gram and tri-gram estimate, the probability of a sentence can be calculated as:

$$P(s) = \prod_{i=1}^n P(w_i/w_{i-1})$$

$$\text{and } P(s) = \prod_{i=1}^n P(w_i/w_{i-2} \dots w_{i-1})$$

As an example, the bi-gram approximation of $P(\text{east}/\text{The Arabian knights are fairy tales of the})$ is

$$P(\text{east}/\text{the}),$$

whereas a tri-gram approximation is

$$P(\text{east}/\text{of the}).$$

A special word (pseudo word) $\langle s \rangle$ is introduced to mark the beginning of the sentence in bi-gram estimation. The probability of the first word in a sentence is conditioned on $\langle s \rangle$. Similarly, in tri-gram estimation, we introduce two pseudo-words $\langle s1 \rangle$ and $\langle s2 \rangle$.

Now, we discuss how to estimate these probabilities. This is done by training the n -gram model on the training corpus. We estimate n -gram parameters using the maximum likelihood estimation (MLE) technique, i.e., using relative frequencies. We count a particular n -gram in the training corpus and divide it by the sum of all n -grams that share the same prefix.

$$P(w_i/w_{i-n+1}, \dots, w_{i-1}) = \frac{C(w_{i-n+1}, \dots, w_{i-1}, w_i)}{\sum_w C(w_{i-n+1}, \dots, w_{i-1}, w)}$$

The sum of all n -grams that share first $n-1$ words is equal to the count of the common prefix $w_{i-n+1}, \dots, w_{i-1}$. So, we rewrite the previous expression as follows:

$$P(w_i/w_{i-n+1}, \dots, w_{i-1}) = \frac{C(w_{i-n+1}, \dots, w_{i-1}, w_i)}{C(w_{i-n+1}, \dots, w_{i-1})}$$

The model parameter we get using these estimates, maximizes the probability of the training set T given the model M , i.e., $P(T|M)$. The frequency with which a word occurs in a text may not be the same as in the training set; this model only provides the most likely solution.

A number of improvements have been suggested for the standard n -gram model. Before we discuss them, let us illustrate these ideas with the help of an example.

Example 2.9

Training set:

The Arabian Knights

These are the fairy tales of the east

The stories of the Arabian knights are translated in many languages

Bi-gram model:

$$\begin{aligned} P(\text{the}/\langle s \rangle) &= 0.67 & P(\text{Arabian}/\text{the}) &= 0.4 & P(\text{knight}/\text{Arabian}) &= 1.0 \\ P(\text{are}/\text{these}) &= 1.0 & P(\text{the}/\text{are}) &= 0.5 & P(\text{fairy}/\text{the}) &= 0.2 \\ P(\text{tales}/\text{fairy}) &= 1.0 & P(\text{of}/\text{tales}) &= 1.0 & P(\text{the}/\text{of}) &= 1.0 \\ P(\text{east}/\text{the}) &= 0.2 & P(\text{stories}/\text{the}) &= 0.2 & P(\text{of}/\text{stories}) &= 1.0 \\ P(\text{are}/\text{knight}) &= 1.0 & P(\text{translated}/\text{are}) &= 0.5 & P(\text{in}/\text{translated}) &= 1.0 \\ P(\text{many}/\text{in}) &= 1.0 & P(\text{languages}/\text{many}) &= 1.0 & \\ P(\text{languages}/\text{many}) &= 1.0 & & & \end{aligned}$$

Test sentence(s): The Arabian knights are the fairy tales of the east.

$$\begin{aligned} P(\text{The}/\langle s \rangle) &\times P(\text{Arabian}/\text{the}) \times P(\text{knight}/\text{Arabian}) \times P(\text{fairy}/\text{the}) \\ &\times P(\text{the}/\text{are}) \times P(\text{tales}/\text{fairy}) \times P(\text{tales}/\text{fairy}) \times P(\text{east}/\text{the}) \\ &= 0.67 \times 0.5 \times 1.0 \times 0.5 \times 0.2 \times 1.0 \times 1.0 \times 1.0 \times 0.2 \\ &= 0.0067 \end{aligned}$$

As each probability is necessarily less than 1, multiplying the probabilities might cause a numerical underflow, particularly in long sentences. To avoid this, calculations are made in log space, where a calculation corresponds to adding log of individual probabilities and taking antilog of the sum.

The n -gram model suffers from data sparseness problem. An n -gram that does not occur in the training data is assigned zero probability, so that even a large corpus has several zero entries in its bi-gram matrix. This is because of the assumption that the probability of occurrence of a word depends only on the preceding word (or preceding $n-1$ words), which is not true in general. There are several long distance dependencies in natural language sentences, which this model fails to capture. Goodman (2003) pointed out that 'there is rarely enough data to accurately estimate the parameters of a language model.'

A number of smoothing techniques have been developed to handle the data sparseness problem, the simplest of these being add-one smoothing. In the words of Jurafsky and Martin (2000):
'Smoothing in general refers to the task of re-evaluating zero probability or low-probability n-grams and assigning them non-zero values.'

The word 'smoothing' is used to denote these techniques because they tend to make distributions more uniform by moving the extreme probabilities towards the average.

2.3.2 Add-one Smoothing

This is the simplest smoothing technique. It adds a value of one to each n -gram frequency before normalizing them into probabilities. Thus, the conditional probability becomes:

$$P(w_i | w_{i-n+1}, \dots, w_{i-1}) = \frac{C(w_{i-n+1}, \dots, w_{i-1}, w_i)}{C(w_{i-n+1}, \dots, w_{i-1}) + V}$$

where V is the vocabulary size, i.e., size of the set of all the words being considered.

In general, add-one smoothing is not considered a good smoothing technique. It assigns the same probability to all missing n -grams, even though some of them could be more intuitively appealing than others. Gale and Church (1994) reported that variance of the counts produced by the add-one smoothing is worse than the unsmoothed MLE method. Another problem with this technique is that it shifts too much of the probability mass towards the unseen n -grams (n -grams with 0 probabilities) as their number is usually quite large. Good-Turing smoothing (Good 1953) attempts to improve the situation by looking at the number of n -grams with a high frequency in order to estimate the probability mass that needs to be assigned to missing or low-frequency n -grams.

SUMMARY

The main topics covered in this chapter are as follows.

- Language modelling deals with providing a description of natural languages amenable to processing.
- There are two main approaches to language modelling: grammar-based language model and statistical language model.
- A grammar-based language model uses grammar to model language.
- A number of computational grammars have been proposed in literature. This chapter provides a discussion of models based on

2.3.3 Good-Turing Smoothing

Good-Turing smoothing (Good 1953) adjusts the frequency f of an n -gram using the count of n -grams having a frequency of occurrence $f+1$. It converts the frequency of an n -gram from f to f^* using the following expression:

$$f^* = (f+1) \frac{n_{f+1}}{n_f}$$

where n_f is the number of n -grams that occur exactly f times in the training corpus. As an example, consider that the number of n -grams that occur 4 times is 25,108 and the number of n -grams that occur 5 times is 20,542. Then, the smoothed count for 5 will be

$$\frac{20542}{25108} \times 5 = 4.09$$

2.3.4 Caching Technique

Another improvement over basic n -gram model is caching. The frequency of n -gram is not uniform across the text segments or corpus. Certain words occur more frequently in certain segments (or documents) and rarely in others. For example, in this section, the frequency of the word 'n-gram' is high, whereas it occurs rarely in earlier sections. The basic n -gram model ignores this sort of variation of n -gram frequency. The cache model combines the most recent n -gram frequency with the standard n -gram model to improve its performance locally. The underlying assumption here is that the recently discovered words are more likely to be repeated.

A number of other smoothing techniques appear in literature. For these, readers are referred to Chen and Goodman (1998).

government and binding, lexical functional grammar, and Paninian grammar.

- GB theory talks about representations at four different levels: structure, d-structure, logical form, and phonetic form.
- An important concept of GB is a general transformational rule called Move α which moves constituents freely, subject to certain constraints (defined by several theories and principles).
- LFG represents sentences at two different levels—constituent structure (*c*-structure) and functional structure (*f*-structure).
- Paninian grammar provides a framework for modelling Indian languages.
- The Paninian framework is syntactico-semantic. The central theme of this framework is Karaka relations.
- A statistical language model estimates the probability (likelihood) of sentence. The most widely used statistical model is n-gram model.
- The n-gram model suffers from sparseness of data. Smoothing techniques such as add-one and Good-Turing can be used to handle this problem.
- Caching is another improvement over the standard n-gram model.

REFERENCES

- Bharati, A., V. Chaitanya, and R. Sangal, 1995, 'Natural Language Processing: A Paninian Perspective', Prentice-Hall of India, New Delhi.
- Bresnan, Joan, 1976, 'Evidence for a theory of unbounded transformation', *Linguistic Analysis*.
- , 1977, 'Variables in theory of transformations', *Formal Syntax*, Peter W. Culicover, Thoman Wasow, and Adrian Akmajian (Eds.), Academic Press, New York, pp. 157–96.
- Chen, Stanley F. and Joshua T. Goodman, 1998, 'An Empirical Study of Smoothing Techniques for Language Modelling', *Technical Report TR-70-98*, Computer Science Group, Harvard University.
- Chomsky, Noam, 1957, *Syntactic Structures*, Mouton, The Hague.
- , 1986, *Some Concepts and Consequences of the Theory of Government and Binding*, MIT Press, Cambridge, MA.
- Gale, William A. and Kenneth W. Church, 1994, 'What's wrong with adding one?', Corpus-based Research into Language, Oostdijk and P. de Haan (Eds.), Rodolpi, Amsterdam.
- Gazdar, G., E. Klein, G. Pullum, and I. Sag, 1985, *Generalized Phrase Structure Grammar*, Blackwell.
- Good, I.J., 1953, 'The population frequencies of species and the estimation of population parameters,' *Biometrika*, 40(3 and 4), pp. 237–64.
- Goodman, Joshua, 2003, 'The state of the art in language modelling,' *Proceedings of the 2003 Conference of the North American Chapter of the Association for Computational Linguistics on Human Language Technology: Tutorials*, 5, Edmonton, Canada.
- Jurafsky, Daniel and James H. Martin, 2000, 'Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics and Speech Recognition', Prentice Hall, New Jersey.
- Kaplan, Ronald M., 1975a, 'On Process Models for Sentence Comprehension', *Explorations in Cognition*, Donald A. Norman and David E. Rumelhart (Eds.), W.H. Freeman, San Francisco.
- , 1975b, 'Transient Processing Load in Relative Clauses', *Doctoral Dissertation*, Harvard University.
- , 1995, 'The Formal Architecture of Lexical-Functional Grammar', *Formal Issues in Lexical-Functional Grammar*, M. Dalrymple, R.M. Kaplan, J.T. Maxwell III, and A. Zaenen (Eds.), CSLI Publications, Stanford.
- Kaplan, Ronald M. and Joan Bresnan, 1982, *The Mental Representation of Grammatical Relations*, Joan Bresnan (Ed.), The MIT Press, Cambridge, MA.
- Kiprasky, P., 1982, 'Some theoretical problem in Panini's grammar', Bhandarkar Oriental Research Institute, Poona.
- Rosenfeld, Ronald, 1994, 'Adaptive Statistical Language Modelling: A Maximum Entropy Approach', *D.Phil. Thesis*, Carnegie Mellon University, Pittsburgh.
- Sells, Peter, 1985, 'Lectures on Contemporary Syntactic Theories', Center for the Study of Language and Information (CSLI), *Lecture Notes*, Number 3, Stanford.
- Shastri, Charudev, 1973, *Vyakarana Chandrodaya*, (Vol. I-V), Motilal Banarsi-dass, in Hindi, New Delhi
- Vasu, S.C. (Tr.), 1977, *The Ashadhyayi of Panini* (2 volumes), Motilal Banarsi-dass, New Delhi.
- Woods, William A., 1970, 'Transition Network Grammars for Natural Language Analysis', *Communications of the ACM*, 13(10), pp. 591–606.
-
1. Give the representation of a sentence in d-structure and s-structure in GB.
 2. How is the structure of a sentence different from the structure of a phrase?

EXERCISES

3. Discuss empty-category principle and give two examples of the creation of empty categories.
4. What is f-structure in LFG? What inputs to an algorithm create an f-structure?
5. Using the annotated rules in Example 2.5, find f-structures for the various phrases and the complete sentence, 'I saw Aparna in the market at night'.
6. What are Karaka relations? Explain the difference between Karta and Agent.
7. What are lexical rules? Give the complete entry for a verb in the lexicon, to be used in LFG.
8. Compare GB and PG. Why is PG called syntactico-semantic theory?
9. What are the problems associated with n-gram model? How are the problems handled?
10. How does caching improve the *n*-gram model?

LAB EXERCISES

1. Create a collection of 10 documents. Write a program to find the frequency of bi-grams in this collection.
2. Find the number of bi-grams that occur more than three times in the collection.