

CHAPTER 3

WORD LEVEL ANALYSIS

CHAPTER OVERVIEW

This chapter focuses on processing carried out at word level, including methods for characterizing word sequences, identifying morphological variants, detecting and correcting misspelled words, and identifying correct part-of-speech of a word. The part-of-speech tagging methods covered in this chapter are: rule-based (linguistic), Stochastic (data-driven), and hybrid.

3.1 INTRODUCTION

As discussed in Chapter 1, natural language processing (NLP) involves different levels and complexities of processing. One way to analyse natural language text is by breaking it down into constituent units (words, phrases, sentences, and paragraphs) and then analyse these units. In Chapter 2, we discussed various language models that are used for analysing the syntax of natural language sentences. Before analysing syntax, we need to understand words, as words are the fundamental unit (syntactic as well as semantic) of any natural language text. This chapter focuses on NLP carried out at word level, including characterizing word sequences, identifying morphological variants, detecting and correcting misspelled words, and identifying correct part-of-speech of a word.

Regular expressions are a beautiful means for describing words. In many text applications, we wish to work with string patterns. Suppose you have just come across the word 'supernova'. It catches your interest and you jump to a search engine to find out more on 'supernovas'. But you do not know whether to type in 'supernova', 'Supernova', or 'supernovas'. Obviously, you need a system, which will retrieve relevant articles using any one of these word forms. Regular expressions are used for describing text strings in situations like this and in other information

retrieval applications. In this chapter, we introduce regular expressions and discuss standard notations for describing text patterns.

After defining regular expressions, we discuss their implementation using finite-state automaton (FSA). Readers who have gone through a course in formal language theory will be familiar with FSA. The FSAs and their variants, such as finite state transducers, have found useful applications in speech recognition and synthesis, spell checking, and information extraction. As we will be using FSA throughout this book, we formally define FSAs in this chapter.

Errors in typing and spelling are quite common in text processing applications. Detecting and correcting these errors are the next topics of discussion in this chapter. Numerous web pages also contain misspelled words and often, query terms entered into the search engines are misspelled. An interactive spelling facility that informs users of such errors and presents appropriate corrections to them, is useful in these applications. There are different classes of words. A word has many forms and the same word may have many different meanings depending on the context. Identifying the class to which a word belongs, its basic form, and its meaning are crucial to analysing text. This is the last topic covered in this chapter.

REGULAR EXPRESSIONS

Regular expressions, or regexes for short, are a pattern-matching standard for string parsing and replacement. They are a powerful way to find and replace strings that take a defined format. For example, regular expressions can be used to parse dates, urls and email addresses, log files, configuration files, command line switches, or programming scripts. They are useful tools for the design of language compilers and have been used in NLP for tokenization, describing lexicons, morphological analysis, etc. We have all used simplified forms of regular expressions, such as the file search patterns used by MS DOS, e.g., dir*.txt.

The use of regular expressions in computer science was made popular by a Unix-based editor, 'ed'. Perl was the first language that provided integrated support for regular expressions. It used a slash around each regular expression; we will follow the same notation in this book. However, slashes are not a part of regular expressions.

Regular expressions were originally studied as a part of theory of computation. They were first introduced by Kleene (1956). A regular expression is an algebraic formula whose value is a pattern consisting of a

set of strings, called the language of the expression. The simplest kind of regular expression contains a single symbol. For example, the expression /a/

denotes the set containing the string 'a'. A regular expression may specify a sequence of characters also. For example, the expression /supernova/

/supernova/

denotes the set that contains the string 'supernova' and nothing else. In a search application, the first instance of each match to regular expression is underlined in Table 3.1.

Table 3.1 Some simple regular expressions

Regular expression	Example patterns
/book/	The world is a book, and those who do not travel read only one page.
/book/	Reporters, who do not read the stylebook, should not criticize their editors.
/face/	Not everything that is faced can be changed. But nothing can be changed until it is faced.
/a/	Region, Observation, and Experience—the Holy Trinity of Science.

3.2.1 Character Classes

Characters are grouped by putting them between square brackets. This way, any character in the class will match one character in the input. For example, the pattern /abcd/ will match (any of) a, b, c, and d. The use of brackets specifies a disjunction of characters. The regular expression [0123456789] specifies any single digit. The character classes are important building blocks in expressions. They sometimes lead to cumbersome notation. For example, it is inconvenient to write the regular expression /abcdefghijklmnopqrstuvwxyz/ to specify 'any lowercase letter'. In these cases, a dash is used to specify a range. The regular expression /[5-9]/ specifies any one of the characters 5, 6, 7, 8, or 9. The pattern /[m-p]/ specifies any one of the letters m, n, o, or p.

Regular expressions can also specify what a single character cannot be, by the use of a caret at the beginning. For example, the pattern /[^x]/ matches any single character except x. This interpretation is true only when a caret appears as the first symbol. If it occurs at any other place, it refers to the caret symbol itself. Table 3.2 shows a few examples explaining these concepts.

Table 3.2 Use of square brackets

RE	Match	Example patterns matched
[abc]	Match any of a, b, and c	'Refresher course will start tomorrow'
[A-Z]	Match any character between A and Z (ASCII order)	'the course will end on Jan. 10, 2006'
[^A-Z]	Match any character other than an uppercase letter	'TREC Conference'
[abc]	Match anything other than a, b, and c	'TREC Conference'
[+*?]	Match any of +, *, ?, or the dot.	'3 ± 2 = 5'
[a^]	Match a or ^	'\n' has three different uses.'

Regular expressions are *case sensitive*. The pattern /s/ matches a lower case 's' but not an uppercase 'S'. This means that the pattern /sana/ will not match the string /Sana/. This problem can be solved by using the disjunction of character s and S. The pattern /[sS]/ will match the strings containing either 's' or 'S'. The pattern /sSana/ matches with the string 'sana' or 'Sana'.

While the use of square brackets solves the capitalization problem, we still need a solution for how to specify both 'supernova' and 'supernovas'.

The pattern /sSupernova[sS]/ matches with any of the strings 'supernovas', 'supernovas', 'Supernovas', and 'SupernovaS', but not with the string 'supernova'. This is achieved with the use of a question mark /?/. A question mark makes the preceding character optional, i.e., zero or one occurrence of the previous character. The regular expression /supernovas/?

specifies both 'supernova' and 'supernovas'. Often, we need to specify repeated occurrences of a character. The * operator, called the *Kleene ** (pronounced 'cleany star'), allow us to do this. The * operator specifies zero or more occurrences of a preceding character or regular expression.

The regular expression /b*/ will match any string containing zero or more occurrences of 'b', i.e., it will match 'b', 'bb', or 'bbbb'. It will also match 'aaa', since that string contains zero occurrences of 'b'. To match a regular expression /bb*/ means a 'b', followed by zero or more 'b's. This will match with any of the strings 'b', 'bb', 'bbb', 'bbbb'. Similarly, the match strings like 'aa', 'bb', or 'abab'. The kleene+ provides a shorter notation to specify one or more occurrences of a character. The regular expression /bb*/ means a 'b', followed by zero or more 'b's. This will match either 'blackberry' or 'blackberries'.

expression /a+/ means one or more occurrences of 'a'. Using Kleene+, we can specify a sequence of digits by the regular expression /[0-9]+/. Complex regular expressions can be built up from simpler ones by means of regular expression operators.

The caret (^) is also used as an anchor to specify a match at the beginning of a line. The dollar sign, \$, is an anchor that is used to specify a match at the end of the line. ^ and \$ are important to regexes. If you wish to search for a line containing only the phrase 'The nature.' and nothing else, you need to specify a regular expression for this search. The anchors ^ and \$ are of great help in this case. The regular expression /The nature\$/ will search exactly for this line.

A number of special characters are also used to build regular expressions. One such character is the dot (.), called wildcard character, which matches any single character. The wildcard expression ./ matches any single character. The regular expression /.at/ matches with any of the string cat, bat, rat, gat, kat, mat, etc. It will also match the meaningless words such as nat, 4at, etc. Table 3.3 shows some of the special characters and their likely use.

Table 3.3 Some special characters

RE	Description
\.	The dot matches any single character.
\n	Matches a new line character (or CR+LF combination).
\t	Matches a tab (ASCII 9).
\d	Matches a digit [0-9].
\D	Matches a non-digit.
\w	Matches an alphanumeric character.
\W	Matches a non-alphanumeric character.
\s	Matches a whitespace character.
\S	Matches a non-whitespace character.
\	Use \ to escape special characters. For example, \. matches a dot, * matches a * and \\ matches a backslash.

We can also use the wildcard symbol for counting characters. For instance /....berry/ matches ten-letter strings that end in berry. This finds patterns like strawberry, sugarberry, and blackberry but fails to match with blueberry or hackberry.

Suppose you are searching a text for the presence of 'Tanveer' or 'Siddiqui'. You cannot use square brackets for this. You need a disjunction operator, shown by the pipe symbol(). The pattern blackberry|blackberries matches either 'blackberry' or 'blackberries'. You might prefer to do this

matching by merely writing blackberry|ies. Unfortunately, this does not work. The pattern blackberry|ies matches with either 'blackberry' or 'ies'. This is because sequences take precedence over disjunction. In order to apply the disjunction operator to a specific pattern, we need to enclose it within parentheses. The parenthesis operator makes it possible to treat the enclosed pattern as a single character for the purposes of neighboring operators. Now, we will consider an example from real application.

Example 3.1 Suppose we need to check if a string is an email address or not. An email address consists of a non-empty sequence of characters followed by the '@' symbol. @ followed by another non-empty sequence of characters ending with pattern like .xx, .xxx, .xxxx, etc. The regular expression for an email address is

$$^*[A-Za-z0-9_.]+@[A-Za-z0-9_.]+[A-Za-z0-9_][A-Za-z0-9_]\$$$

Table 3.4 shows the various parts of this 'rgex'.

Table 3.4 Parts of regular expression of Example 3.1

Pattern	Description
^*[A-Za-z0-9_.]+	Match a positive number of acceptable characters at the start of the string.
@	Match the @ sign.
[A-Za-z0-9_.]+	Match any domain name, including a dot.
[A-Za-z0-9_.]+\\$	Match two acceptable characters but not a dot. This ensures that the email address ends with .xx, .xxx, .xxxx, etc.

This example works for most cases. However, the specification is not based on any standard and may not be accurate enough to match all correct email addresses. It may accept non-working email addresses and reject working ones. Fine-tuning is required for accurate characterization.

A regular expression characterizes a particular kind of formal language, called a regular language. The language of regular expressions is similar to formulas of Boolean logic. Like logic formulas, regular expressions represent sets. Regular language is set of strings described by the regular expression. Regular languages may be encoded as finite state networks.

A regular expression may contain symbol pairs. For example, the regular expression /a:b/ represents a pair of strings. The regular expression /a:b/ actually denotes a regular relation. A regular relation may be viewed as a mapping between two regular languages. The a:b relation is simply the cross product of the languages denoted by the expressions /a/ and /b/. To differentiate the two languages that are involved in a regular relation, we call the first one, the upper, and the second one, the lower

language, of the relation. Similarly, in the pair /a:b/, the first symbol, a, can be called the upper symbol and the second symbol, b, the lower symbol. The two components of a symbol pair are separated in our notation by a colon (:). Without any whitespace before or after. To make the notation less cumbersome, we ignore the distinction between the language A and the identity relation that maps every string of A to itself. Therefore, we also write /a:a/ simply as /a/.

Regular expressions have clean, declarative semantics (Voutilainen 1996). Mathematically, they are equivalent to finite automata, both having the same expressive power. This makes it possible to encode regular languages using finite-state automata, leading to easier manipulation of context free and other complex languages. Similarly, regular relations can be represented using finite-state transducers. With this representation, it is possible to derive new regular languages and relations by applying regular operators, instead of re-writing the grammars.

3.3 FINITE-STATE AUTOMATA

In our childhood, each of us must have played some game that fits the following description:

Pieces are set up on a playing board; dice are thrown or a wheel is spun, and a number is generated at random. Based on the number appearing on the dice, the pieces on the board are rearranged specified by the rules of the game. Then, it is your opponent's turn; she also does the same thing, resulting in rearrangement of the pieces based on the number generated. There is no skill or choice involved. The entire game is based on the values of the random numbers.

Consider all possible positions of the pieces on the board and call them states. The state in which the game begins is termed the *initial state*, and the state corresponding to the winning positions is termed the *final state*. We begin with the initial state of the starting positions of the pieces on the board. The game then changes from one state to another based on the value of the random number. For each possible number, there is one and only one resulting state given the input of the number and the current state. This continues until one player wins and the game is over. This is called a final state.

Now consider a very simple machine with an input device, a processor, some memory, and an output device. The machine starts in the initial state. It checks the input and goes to next state, which is completely determined by the prior state and the input. If all goes well, the machine

reaches final state and terminates. If the machine gets an input for which the next state is not specified, and it gets stuck at a non-final state, we say the machine has failed or rejected the input.

A general model of this type of machine is called a finite automaton; 'finite', because the number of states and the alphabet of input symbols is finite; 'automaton' because the machine moves automatically, i.e., the change of state is completely governed by the input. This type of machine is more commonly called deterministic.

A finite automaton has the following properties:

1. A finite set of states, one of which is designated the initial or start state, and one or more of which are designated as the final states.
2. A finite alphabet set, Σ , consisting of input symbols.
3. A finite set of transitions that specify for each state and each symbol of the input alphabet, the state to which it next goes.

A finite automaton can be deterministic or non-deterministic. In a non-deterministic automaton, more than one transition out of a state is possible for the same input symbol.

Example 3.2 Suppose $\Sigma = \{a, b\}$, the set of states $= \{q_0, q_1, q_2, q_3, q_4\}$ with q_0 being the start state and q_4 the final state, we have the following rules of transition:

1. From state q_0 and with input a , go to state q_1 .
2. From state q_1 and with input b , go to state q_2 .
3. From state q_1 and with input c , go to state q_3 .
4. From state q_2 and with input b , go to state q_4 .
5. From state q_3 and with input b , go to state q_4 .

This finite-state automaton is shown as a directed graph, called transition diagram, in Figure 3.1. The nodes in this diagram correspond to the final states, and the arcs to transitions. The arcs are labelled with inputs. The final state is represented by a double circle. As seen in the figure, there is exactly one transition leading out of each state. Hence, this automaton is deterministic.

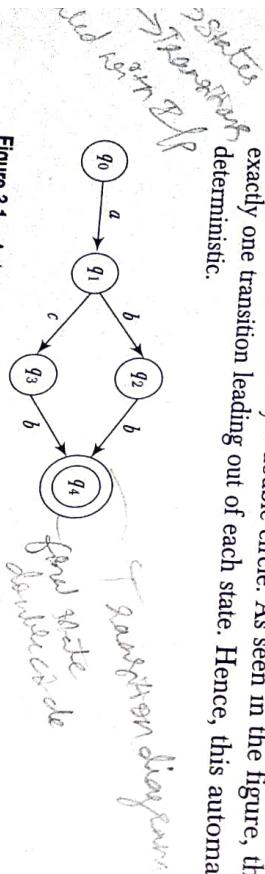


Figure 3.1 A deterministic finite-state automaton (DFA)

Finite-state automata have been used in a wide variety of areas, including linguistics, electrical engineering, computer science, mathematics, and logic. These are an important tool in computational linguistics and have been used as a mathematical device to implement regular expressions. Any regular expression can be represented by a finite automaton and the language of any finite automaton can be described by a regular expression. Both have the same expressive power. The following formal definitions of the two types of finite state automaton, namely, deterministic and non-deterministic finite automaton, are taken from Hopcroft and Ullman (1979).

A deterministic finite-state automaton (DFA) is defined as a 5-tuple $(Q, \Sigma, \delta, S, F)$, where Q is a set of states, Σ is an alphabet, S is the start state, $F \subseteq Q$ is a set of final states, and δ is a transition function. The transition function δ defines mapping from $Q \times \Sigma$ to Q . That is, for each state q and symbol a , there is at most one transition possible as shown in Figure 3.1.

Unlike DFA, the transition function of a non-deterministic finite-state automaton (NFA) maps $Q \times (\Sigma \cup \{\epsilon\})$ to a subset of the power set of Q . That is, for each state, there can be more than one transition on a given symbol, each leading to a different state.

This is shown in Figure 3.2, where there are two possible transitions from state q_0 on input symbol a .

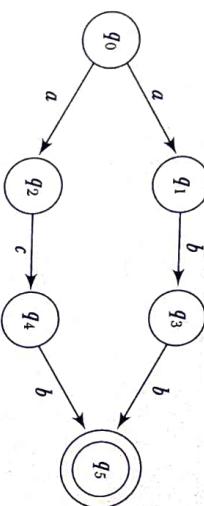


Figure 3.2 Non-deterministic finite-state automaton (NFA)

A path is a sequence of transitions beginning with the start state. A path leading to one of the final states is a successful path. The FSAs encode regular languages. The language that an FSA encodes is the set of strings that can be formed by concatenating the symbols along each successful path. Clearly, for automata with cycles, these sets are not finite.

We now examine what happens to various input strings that are presented to finite state automata. Consider the deterministic automaton described in Example 3.2 and the input ac. We start with state q_0 and go to state q_1 . The next input symbol is c , so we go to state q_3 . No more input is left and we have not reached the final state, i.e., we have an unsuccessful end. Hence, the string ac is not recognized by the automaton.

This example illustrates how an FSA can be used to accept or recognize a string. The set of all strings that leave us in a final state is called the language accepted or defined by the FA. This means *ac* is not a word in the language defined by the automaton of Figure 3.1.

Now, consider the input acb . Again, we start with state q_0 and go to state q_1 . The next input symbol is c , so we go to state q_3 . The next input symbol is b , which leads to state q_4 . No more input is left and we have reached to final state, i.e., this time we have a successful termination. Hence, the string abc is a word of the language defined by the automaton.

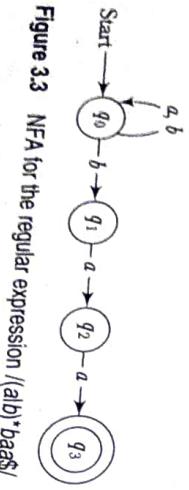
The language defined by this automaton can be described by the regular expression $/ab\bar{b}ach/$.

The example considered here is quite simple. Typically, the list of transition rules can be quite long. Listing all transition rules may be inconvenient, so often we represent an automaton as a *state-transition table*. The rows in this table represent states and the columns correspond to input. The entries in the table represent the transition corresponding to a given state-input pair. A ϕ entry indicates missing transition. This table contains all the information needed by FSA. The state transition table for the automaton considered in Example 3.2 is shown in Table 3.5.

Table 3.5 The state-transition table for the DFA shown in Figure 3.11

State	a	b	c	Input
start:	q_0	q_1	ϕ	
	q_1	ϕ	q_2	q_3
	q_2	ϕ	q_4	ϕ
	q_3	ϕ	q_4	ϕ
final:	q_4	ϕ	ϕ	

language consisting of all strings containing only *a*s and *b*s and ending with *baa*. We can specify this language by the regular expression $(ab)^*baa\$$. The NFA implementing this regular expression is shown in Figure 3.3.



• **QUESTION** What is the regular expression `(\alh)*baagf`?

There are two broad classes of morphemes: stems and affixes. The stem is the main morpheme, i.e., the morpheme that contains the central meaning. Affixes modify the meaning given by the stem. Affixes are divided into prefix, suffix, infix, and circumfix. Prefixes are morphemes which appear before a stem, and suffixes are morphemes applied to the end of the stem. Circumfixes are morphemes that may be applied to either end of the stem while infixes are morphemes that appear inside a stem. Prefixes and suffixes are quite common in Urdu, Hindi, and English. For example, the Urdu word بُرْت (beauty), meaning untimely, is composed of the stem, wag, and the prefix, be, بے. گھوڑا (ghodha) is composed of the stem, گھوڑا (ghodha) and the suffix, اں (on). Also, the word زارووت مند (zarooratmand), is composed of the stem, زارووت (zaroorat) and the prefix, مند (mand).

3.4 MORPHOLOGICAL PARSING

Morphology is a sub-discipline of linguistics. It studies word structure and the formation of words from smaller units (morphemes). The goal of morphological parsing is to discover the morphemes that build a given word. Morphemes are the smallest meaning-bearing units in a language. For example, the word 'bread' consists of a single morpheme and 'eggs' consist of two: the morpheme egg and the morpheme-s. A morphological parser should be able to tell us that the word 'eggs' is the plural form of

There are two broad classes of morphemes: stems and affixes. The stem is the main morpheme, i.e., the morpheme that contains the central meaning. Affixes modify the meaning given by the stem. Affixes are divided into prefix, suffix, infix, and circumfix. Prefixes are morphemes which appear before a stem, and suffixes are morphemes applied to the end of the stem. Circumfixes are morphemes that may be applied to either end of the stem while infixes are morphemes that appear inside a stem. Prefixes and suffixes are quite common in Urdu, Hindi, and English.

For example, the Urdu word, بُرُوت (buroot), meaning untimely, is composed of the stem, *waqt*, and the prefix, *be-*. بُرُوت (buroot) is composed of the stem, گھوڑا (*ghodha*) and the suffix, وں (*on*). Also, the word, ضرورتمند (zaroortmand), is composed of the stem, ضرورت (zaroort).

Table 3.6 The state transition matrix

State	Input	
	a	b
start:		
q_0	$\{q_0\}$	$\{q_0, q_1\}$
q_1	$\{q_2\}$	ϕ
q_2	$\{q_3\}$	ϕ
q_3	ϕ	ϕ
final:		

An NFA can be converted to an equivalent DFA and vice versa. The equivalent DFA for the NFA shown in Figure 3.3 is shown in Figure 3.4.

(*unhappy*), and the suffix, **मंड** (*mand*). Similarly, the English word, *unhappy*, is composed of the stem, *happy*, and the prefix, *un-*. The English word, *bird*, is composed of the stem, *bird*, and the suffix, *-s*. Likewise, the Hindi word, **गिरिजा** is composed of a stem **गिरि** and the suffix – **जा**.

Telugu word **గुర్రములు** (*gurramulu*-plural form of *gurramu*, meaning 'horse') is composed of the stem **గుర్రము** (*gurramu*) and suffix -**లు** (*lu*). Some commonly used suffixes in plural forms of Telugu nouns are: **లు** (*lu*), **లూ** (*luu*), **లౌ** (*luu*). Here is a list of singular and plural forms of a few Telugu words:

Singular	Meaning	Plural
పిల్ల (pilli)	Cat	పిల్లలు (pillulu)
పదుచు (paduchu)	Women	పదుచులు (paduchulu)
ఆడది (aadadi)	Women	ఆడవాండు (aadavandu* or aadivalu)
మగానుడు (maganudu)	Man	మగావండు (magavandu or magavallu)
పురుషుడు (purushudu)	Man	పురుషులు (purushulu)
చెఱి (chevu)	Ear	చెఱులు (chevulu)
ఇల్లు (illu)	House	ఇల్లలు (illulu or illu)

*Another plural form of *aadadi* is *aadavaru*, which is used to show respect.

There are three main ways of word formation: inflection, derivation, and compounding. In inflection, a root word is combined with a grammatical morpheme to yield a word of the same class as the original stem. Derivation combines a word stem with a grammatical morpheme to yield a word belonging to a different class, e.g., formation of the noun 'computation' from the verb 'compute'. The formation of a noun from a verb or adjective is called nominalization. Compounding is the process of merging two or more words to form a new word. For example, personal computer, desktop, overlook. Morphological analysis and generation deal with inflection, derivation and compounding process in word formation. New words are continually forming a natural language. Many of these are morphologically related to known words. Understanding morphology is therefore important to understand the syntactic and semantic properties of new words. Morphological analysis and generation are essential to many NLP applications ranging from spelling corrections to machine translations. In parsing, e.g., it helps to know the agreement features of words. In information retrieval, morphological analysis helps identify the presence of a query word in a document in spite of different morphological variants.

Parsing, in general, means taking an input and producing some sort of structures for it. In NLP, this structure might be morphological, syntactic, semantic, or pragmatic. Morphological Parsing takes as input the inflected

Inflected Surface Form of Each Word

surface form of each word in a text. As output, it produces the parsed form consisting of a canonical form (or *lemma*) of the word and a set of tags showing its syntactical category and morphological characteristics, e.g., possible part of speech and/or inflectional properties (gender, number, person, tense, etc.). Morphological generation is the inverse of this process.

Both analysis and generation rely on two sources of information: a dictionary of the valid lemmas of the language and a set of inflection paradigms.

A morphological parser uses following information sources:

1. Lexicon

A lexicon lists stems and affixes together with basic information about them.

2. Morphotactics

There exists certain ordering among the morphemes that constitute a word. They cannot be arranged arbitrarily. For example, rest-lessness is a valid word in English but not rest-less. Morphotactics deals with the ordering of morphemes. It describes the way morphemes are arranged or touch each other.

3. Orthographic rules

These are spelling rules that specify the changes that occur when two given morphemes combine. For example the *y* → *ier* spelling rule changes 'easy' to 'easier' and not to 'eayer'.

Morphological analysis can be avoided if an exhaustive lexicon is available that lists features for all the word-forms of all the roots. Given a word, we simply consult the lexicon to get its feature values. For example, suppose an exhaustive lexicon for Hindi contains the following entries for the Hindi root-word *ghoonthaa*.

Table 3.7 A sample lexicon entry

Word form	Category	Root	Gender	Number	Person
Ghoonthaa	noun	GhoDaa	masculine	singular	3rd
Ghoonthii	-do-	-do-	feminine	-do-	-do-
Ghoonthon	-do-	-do-	masculine	plural	-do-
Ghoonthee	-do-	-do-	-do-	-do-	-do-

Given a word, say *ghoonthon*, we can look up the lexicon to get its feature values.

However, this approach has several limitations. First, it puts a heavy demand on memory. We have to list every form of the word, which results in a large number of, often redundant, entries in the lexicon.

Second, an exhaustive lexicon fails to show the relationship between different roots having similar word-forms. That means the approach fails to capture linguistic generalization, which is essential to develop a system capable of understanding unknown words.

Third, for morphologically complex languages, like Turkish, the number of possible wordforms may be theoretically infinite. It is not practical to list all possible word forms in these languages.

These limitations explain why morphological parsing is necessary. The complexity of the morphological analysis varies widely among the world's languages, and is quite high even in relatively simple cases, such as English.

The simplest morphological systems are stemmers that collapse morphological variations of a given word (word-forms) to one lemma or stem. They do not require a lexicon. Stemmers have been especially used in information retrieval. Two widely used stemming algorithms have been developed by Lovins (1968) and Porter (1980). Stemmers do not use a lexicon; instead, they make use of rewrite rules of the form:

$$ier \rightarrow y \text{ (e.g., earlier} \rightarrow \text{early)}$$

$$ing \rightarrow \epsilon \text{ (e.g., playing} \rightarrow \text{play)}$$

Stemming algorithms work in two steps:

- (i) Suffix removal: This step removes predefined endings from words.

- (ii) Recoding: This step adds predefined endings to the output of the first step.

These two steps can be performed sequentially as in Lovins's stemmer or simultaneously as in Porter's stemmer. For example, Porter's stemmer makes use of the following transformation rule:

$$ational \rightarrow ate$$

to transform word such as 'rotational' into 'rotate'.

It is difficult to use stemming with morphologically rich languages. Even in English, stemmers are not perfect. Krovitz (1993) pointed out errors of omissions and commissions in the Porter algorithm, such as 'noisy'. Another problem with Porter's algorithm is that it reduces only suffixes, prefixes and compounds are not reduced.

A more efficient two-level morphological model, first proposed by Koskenniemi (1983), can be used for highly inflected languages. In this model, a word is represented as a correspondence between its lexical level form and its surface level form. The surface level represents the actual spelling of the word while the lexical level represents the concatenation of its constituent morphemes. Morphological parsing is viewed as a mapping from the surface level into morpheme and feature sequences on the lexical level.

Two-level
morphological
model

Surface level	p	l	a	y	i	n	g
---------------	---	---	---	---	---	---	---

Lexical level	p	l	a	y	+V	PP
---------------	---	---	---	---	----	----

Figure 3.5 Surface and lexical forms of a word

Similarly, the surface form 'books' is represented in the lexical form as 'book + N + PL', where the first component is the stem and the second component (N + PL) is the morphological information, which tells us that the surface level form is a plural noun.

This model is usually implemented with a kind of finite-state automata, called finite-state transducer (FST). A transducer maps a set of symbols to another. A finite state transducer does this through a finite state automaton. An FST can be thought of as a two-state automaton, which recognizes or generates a pair of strings. FST passes over the input string by consuming the input symbols on the tape it traverses and consists it to the output string in the form of symbols. Formally, an FST has been defined by Hopcroft and Ullman (1979) as follows:

A finite-state transducer is a 6-tuple $(\Sigma_1, \Sigma_2, Q, \delta, S, F)$, where θ is set of states, S is the initial state, and $F \subseteq Q$ is a set of final states, Σ_1 is input alphabet, Σ_2 is output alphabet, and δ is a function mapping $Q \times (\Sigma_1 \cup \{\epsilon\}) \times (\Sigma_2 \cup \{\epsilon\})$ to a subset of the power set of Q .

Transducers can be seen as automata with transitions labelled with symbols from $\Sigma_1 \times \Sigma_2$, where Σ_1 and Σ_2 are the alphabets of input and output respectively. Thus, an FST is similar to an NFA except in that transitions are made on strings rather than on symbols and, in addition, they have outputs.

Figure 3.6 shows a simple transducer that accepts two input strings, hot and cat, and maps them onto cot and bat respectively. It is a common practice to represent a pair like $a:a$ by a single letter.



Figure 3.6 Finite-state transducer

Model Language Just as FSAs encode regular languages, FSTs encode regular relations. Regular relation is the relation between regular languages. The regular language encoded on the upper side of an FST is called upper language. If T is a relation and the one on the lower side is termed lower language. If T is a transducer, and s is a string, then we use $T(s)$ to represent the set of strings encoded by T such that the pair (s, t) is in the relation.

The FSTs are closed under union concatenation, composition, and Kleene closure. However, in general, they are not closed under intersection and complementation.

With this introduction, we can now implement the two-level morphology using FST. To get from the surface form of a word to its morphological analysis, we proceed in two steps as illustrated in Figure 3.7.

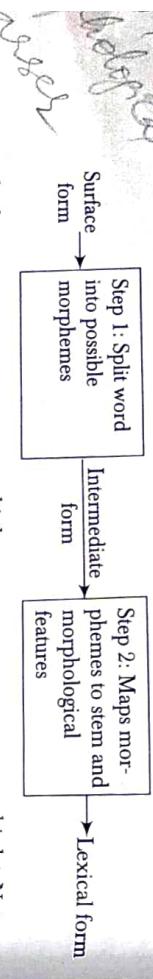


Figure 3.7 Two-step morphological parser

First, we split the words up into its possible components. For example, we make *bird* + *s* out of *birds*, where + indicates morpheme boundaries.

In this step, we also consider spelling rules. Thus, there are two possible ways of splitting up *boxes*, namely *boxe* + *s* and *box* + *s*. The first one assumes that *boxe* is the stem and *s* the suffix, while the second assumes that *box* the stem is and that *e* has been introduced due to the spelling rule. The output of this step is a concatenation of morphemes, i.e., stems and affixes. There can be more than one representation for a given word. A transducer that does the mapping (translation) required by this step for the surface form '*lesser*' might look like Figure 3.8. This FST represents the information that the comparative form of the adjective *less* is *lesser*; *e* here is the empty string. The automaton is inherently bi-directional: the same transducer can be used for analysis (surface input, 'upward' application) or for generation (lexical input, 'downward' application).

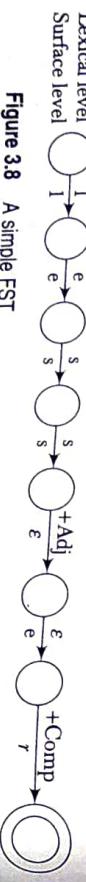


Figure 3.8 A simple FST

In the second step, we use a lexicon to look up categories of the stems and meaning of the affixes. So, *bird* + *s* will be mapped to *bird+N+PL*, and *box* + *s* to *box+N+PL*. We also find out now that *boxe* is not a legal stem. This tells us that splitting *boxes* into *boxe* + *s* is an incorrect way of splitting *boxes*, and should therefore be discarded. This may not be the case always. We have words like *spouse* or *parses* where splitting the word into *spouse* + *s* or *parses* + *s* is correct. Orthographic rules are used to handle these spelling variations. For instance, one of the spelling rules says add *e* after *-s*, *-z*, *-x*, *-ch*, *-sh* before the *s* (e.g., *dish* → *dishes*, *box* → *boxes*). Each of these steps can be implemented with the help of a transducer. Thus, we need to build two transducers: one that maps the surface form to the intermediate form and another that maps the intermediate form to the lexical form.

We now develop an FST-based morphological parser for singular and plural nouns in English. The plural form of regular nouns usually end with *-s* or *-es*. However, a word ending in '*s*' need not necessarily be the plural form of a word. There are a number of singular words ending in *s*, e.g., *miss* and *ass*. One of the required translations is the deletion of the '*e*' when introducing a morpheme boundary. This deletion is usually required for words ending in *xes*, *ses*, *zes* (e.g., suffixes and *boxes*). The transducer in Figure 3.9 does this. Figure 3.10 shows the possible sequences of states that the transducer undergoes, given the surface forms *birds* and *boxes* as input.

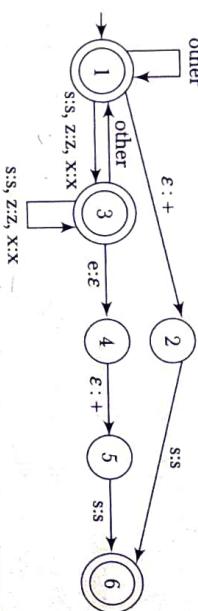


Figure 3.9 A simplified FST, mapping English nouns to the intermediate form

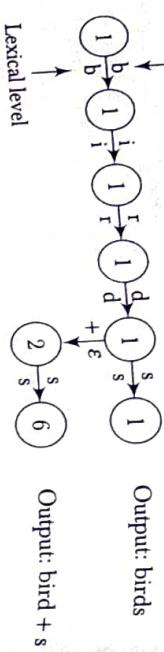
The next step is to develop a transducer that does the mapping from the intermediate level to the lexical level. The input to transducer has one of the following forms:

- Regular noun stem, e.g., *bird*, *cat*
- Regular noun stem + *s*, e.g., *bird+s*
- Singular irregular noun stem, e.g., *goose*
- Plural irregular noun stem, e.g., *geese*

In the first case, the transducer has to map all symbols of the stem to themselves and then output *N* and *g* (Figure 3.7). In the second case, it



Figure 3.10 Possible sequences of states



Input: boxes

Output: box + s

Input: birds

Output: birds

Lexical level

Surface level

Lexical level

Surface level

Mapping for the regular surface form of *bird* is $b:b\ i:i\ r:r\ d:d$. Representing pairs like *a:a* with a single letter, these two representations are reduced to $g\ e:o\ e:o\ s:s$ and $b\ i:r\ d$ respectively.

Composing this transducer with the previous one, we get a single two-level transducer with one input tape and one output tape. This maps plural nouns into the stem plus the morphological marker $+pl$ and singular nouns into the stem plus the morpheme $+sg$. Thus a surface word form *birds* will be mapped to *bird + N + pl* as follows:

$$b:b\ i:i\ r:r\ d:d + \epsilon:N + s:pl$$

Each letter maps to itself, while ϵ maps to morphological feature $+N$, and s maps to morphological feature pl . Figure 3.12 shows the resulting composed transducer.

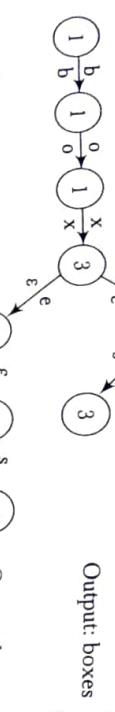


Figure 3.12 A transducer mapping nouns to their stem and morphological features

The power of the transducer lies in the fact that the same transducer can be used for analysis and generation. That is, we can run it in the downward direction (input: surface form and output: lexical form) or in the upward direction.

Figure 3.11 Transducer for Step 2

has to map all symbols of the stem to themselves, but then output N and

replaces PL with s . In the third case, it has to do the same as in the first case. Finally, in the fourth case, the transducer has to map the irregular plural noun stem to the corresponding singular stem (e.g., *geese* to *goose*) and then it should add N and PL . The general structure of this transducer looks like Figure 3.11.

The mapping from State 1 to State 2, 3, or 4 is carried out with the help of a transducer encoding a lexicon. The transducer implementing the lexicon maps the individual regular and irregular noun stems to their correct noun stem, replacing labels like regular noun form, etc. This lexicon maps the surface form *goose*, which is an irregular noun, to its correct stem *goose* in the following way:

$$gg\ e:o\ e:o\ s:s\ e:e$$

3.5 SPELLING ERROR DETECTION AND CORRECTION

In computer-based information systems, errors of typing and spelling constitute a very common source of variation between strings. These errors have been widely investigated. All investigations agree that single character omission, insertion, substitution, and reversal are the most common typing mistakes. In an early investigation, Damearnu (1964) reported that over 80% of the typing errors were single-error misspellings:

- (1) substitution of a single letter,
- (2) omission of a single letter,
- (3) insertion of a single letter, and
- (4) transposition of two adjacent letters.

Shafer and Hardwick (1968) found that the most common type of single character error was substitution, followed by omission of a letter, and then insertion of a letter. Single character omission occurs when a single character is missed (deleted), e.g., when 'concept' is accidentally typed as 'concpt'. Insertion error refers to the presence of an extra character in a word, e.g. when 'error' is misspell as 'errorn'. Substitution error occurs when a wrong letter is typed in place of the right one, as in 'errpr', where 'p' appears in place of 'o'. Reversal refers to a situation in which the sequence of characters is reversed, e.g., 'aer' instead of 'are'. This is also termed transposition.

Optical character recognition (OCR) and other automatic reading devices introduce errors of substitution, deletion, and insertion but not of reversal. OCR errors are usually grouped into five classes: substitution, multi-substitution (or framing), space deletion or insertion, and failures. Unlike substitution errors made in typing, OCR substitution errors are caused due to visual similarity such as c→e, l→l, r→n. The same is true for multi-substitution, e.g., m→rn. Failure occurs when the OCR algorithm fails to select a letter with sufficient accuracy. The frequency and type of errors are characteristics of the particular device. These errors can be corrected using 'context' or by using linguistic structures.

Many approaches to speech recognition deal with strings of phonemes (or symbols representing sounds), and attempt to match a spoken utterance with a dictionary of known utterances.

Unlike typing errors, spelling errors are mainly phonetic, where the misspell word is pronounced in the same way as the correct word. Phonetic errors are harder to set right because they distort the word by more than a single insertion, deletion, or substitution. Phonetic variations are common in transliteration. For example,

Spelling errors belong to one of two distinct categories: non-word errors and real word errors. When an error results in a word that does not appear in a given Lexicon or is not a valid orthographic word form, it is termed a non-word error. Most of the early research on spelling errors focused on the detection of such non-words. The two main techniques used were n-gram analysis and dictionary lookup. Non-word error detection is now considered a solved problem.

A real-word error results in actual words of the language. It occurs due to typographical mistakes, or spelling errors, e.g., substituting the spelling of a homophone or near-homophone, such as piece for peace or meat for meet. Real-word errors may cause local syntactic errors, global syntactic

30

errors, semantic errors, or errors at discourse or pragmatic levels. It becomes impossible to decide that a word is wrong without some contextual information.

Spelling correction consists of detecting and correcting errors. Error detection is the process of finding misspelled words and error correction is the process of suggesting correct words to a misspelled one. These sub-problems are addressed in two ways:

1. Isolated-error detection and correction
2. Context-dependent error detection and correction

In isolated-word error detection and correction, each word is checked separately, independent of its context. Detecting whether or not a word is correct seems simple—why not to look up the word in a lexicon? Unfortunately, it is not as simple as it appears. There are a number of problems associated with this simple strategy.

- The strategy requires the existence of a lexicon containing all correct words. Such a lexicon would take a long time to compile and occupy a lot of space.
- Some languages are highly productive. It is impossible to list all the correct words of such languages.

• This strategy fails when spelling error produces a word that belongs to the lexicon, e.g., when 'theses' is written in place of 'these'. Such an error is called a *real-word error*.

- The larger the lexicon, the more likely it is that an error goes undetected, because the chance of a word being found is greater in a large lexicon.

Context dependent error detection and correction methods, utilize the context of a word to detect and correct errors. This requires grammatical analysis and is thus more complex and language dependent. Even in context dependent methods, the list of candidate words must first be obtained using an isolated-word method before making a selection depending on the context.

The spelling correction algorithm has been broadly categorized by Kukich (1992) as follows.

Minimum edit distance The minimum edit distance between two strings is the minimum number of operations (insertions, deletions, or substitutions) required to transform one string into another. Spelling correction algorithms based on minimum edit distance are the most studied algorithms.

Similarity key techniques The basic idea in a similarity key technique is to change a given string into a key such that similar strings will change into the same key. The SOUNDEX system (Odell and Russell 1918) is an example of a system that uses this technique in phonetic spelling correction applications.

n-gram based techniques The n-grams can be used for both non-word and real-word error detection because in the English alphabet, certain bi-grams and tri-grams of letters never occur or rarely do so; for example the tri-gram *qyt* and the bi-gram *qd*. This information can be used to handle non-word error. Strings that contain these unusual n-grams can be identified as possible spelling errors. n-gram techniques usually require a large corpus or dictionary as training data, so that an n-gram table of possible combinations of letters can be compiled. In case of real-word error detection, we calculate the likelihood of one character following another and use this information to find possible correct word candidates.

Neural nets These have the ability to do associative recall based on incomplete and noisy data. They can be trained to adapt to specific spelling error patterns. The drawback of neural nets is that they are computationally expensive.

Rule-based techniques In a rule-based technique, a set of rules (heuristics) derived from knowledge of a common spelling error pattern is used to transform misspelled words into valid words. For example, if it is known that many errors occur from the letters *ue* being typed as *eu*, then we may write a rule that represents this.

3.5.1 Minimum Edit Distance

The minimum edit distance is the number of insertions, deletions, and substitutions required to change one string into another (Wagner and Fischer 1974). When we talk about distance between two strings, we are talking of the minimum edit distance. For example, the minimum edit distance between 'tutor' and 'tumour' is 2: We substitute 'm' for 't' and conversion. Therefore, the minimum edit sequence can be found for this

between two strings can be represented as a binary function, *ed*, which maps two strings to their edit distance. *ed* is symmetric. For any two strings, *s* and *t*, *ed(s, t)* is always equal to *ed(t, s)*.

Edit distance can be viewed as a string alignment problem. By aligning two strings, we can measure the degree to which they match. There may be more than one possible alignment between two strings. The best

possible alignment corresponds to the minimum edit distance between the strings. The alignment shown here, between *tutor* and *tumour*, has a distance of 2.

A dash in the upper string indicates insertion. A substitution occurs when the two alignment symbols do not match (shown in bold). We can associate a weight or cost with each operation. The Levenshtein distance between two sequences is obtained by assigning a unit cost to each operation. Another possible alignment for this sequences is:

t	u	t	o	-	r
t	u	-	m	o	u

diff copy

which has a cost of 3. We already have a better alignment than this one.

The problem of finding minimum edit distance seems quite simple but in fact is not so. A choice that seems good initially might lead to problems later. Dynamic programming algorithms can be quite useful for finding minimum edit distance between two sequences. Dynamic programming refers to a class of algorithms that apply a table-driven approach to solve problems by combining solutions to sub-problems. The dynamic programming algorithm for minimum edit distance is implemented by creating an edit distance matrix. This matrix has one row for each symbol in the source string and one column for each symbol in the target string. The (i, j) th cell in this matrix represents the distance between the first *i* character of the source and the first *j* character of the target string. Each cell can be computed as a simple function of its surrounding cells. Thus, by starting at the beginning of the matrix, it is possible to fill each entry iteratively. The value in each cell is computed in terms of three possible paths.

$$dist[i, j] = \begin{cases} dist[i - 1, j] + insert_cost, \\ dist[i - 1, j - 1] + subst_cost[source_i, target_j], \\ dist[i, j - 1] + delete_cost \end{cases}$$

The substitution will be 0 if the *i*th character in the source matches with *j*th character in the target. The minimum edit distance algorithm is shown in Figure 3.13. How the algorithm computes the minimum edit distance between *tutor* and *tumour* is shown in Figure 3.14.

```

Input: Two strings,  $X$  and  $Y$ 
Output: The minimum edit distance between  $X$  and  $Y$ 
 $m \leftarrow \text{length}(X)$ 
 $n \leftarrow \text{length}(Y)$ 
for  $i = 0$  to  $m$  do
     $\text{dist}[i, 0] \leftarrow i$ 
for  $j = 0$  to  $n$  do
     $\text{dist}[0, j] \leftarrow j$ 
for  $i = 0$  to  $m$  do
    for  $j = 0$  to  $n$  do
         $\text{dist}[i, j] = \min\{ \text{dist}[i-1, j] + \text{insert\_cost},$ 
         $\text{dist}[i-1, j-1] + \text{subst\_cost}(X_i, Y_j),$ 
         $\text{dist}[i, j-1] + \text{delete\_cost} \}$ 

```

Figure 3.13 Minimum edit distance algorithm

#	t	u	m	o	u	r
0	1	2	3	4	5	6
1	0	1	2	3	4	5
2	1	0	1	2	3	4
3	2	1	1	2	3	4
4	3	2	2	1	2	3
5	4	3	3	2	2	2
6	5	4	3	3	2	2

Figure 3.14 Computing minimum edit distance

Minimum edit distance algorithms are also useful for determining accuracy in speech recognition systems. Kemal Oflazer (1996) proposed an efficient algorithm based on spelling correction with finite-state automata.

6 WORDS AND WORD CLASSES

Words are classified into categories called part-of-speech. These are sometimes called word classes or lexical categories. These lexical categories are usually defined by their syntactic and morphological behaviours. The most common lexical categories are nouns and verbs. Other lexical categories include adjectives, adverbs, prepositions, and conjunctions. Table 3.8 shows some of the word classes in English. Lexical categories and their properties vary from language to language. Word classes are further categorized as open and closed word classes. Open word classes constantly

Part-of-speech example	noun	NN	student, chair, proof, mechanism
VB	verb	VB	study, increase, produce
ADJ	adj	ADJ	large, high, tall, few,
JJ	adverb	JJ	carefully, slowly, uniformly
IN	preposition	IN	in, on, to, of
PRP	pronoun	PRP	I, me, they
DET	determiner	DET	the, a, an, this, those

3.7 PART-OF-SPEECH TAGGING

Part-of-speech tagging is the process of assigning a part-of-speech (such as a noun, verb, pronoun, preposition, adverb, and adjective), to each word in a sentence. The input to a tagging algorithm is the sequence of words of a natural language sentence and specified tag sets (a finite list of part-of-speech tags). The output is a single best part-of-speech tag for each word. Many words may belong to more than one lexical category. For example, the English word ‘book’ can be a noun as in ‘I am reading a good book’ or a verb as in ‘The police booked the snatcher’. The same is true for other languages. For example, the Hindi word ‘soma’ may mean ‘gold’ (noun) or ‘sleep’ (verb). However, only one of the possible meanings is used at a time. In tagging, we try to determine the correct lexical category of a word in its context. No tagger is efficient enough to identify the correct lexical category of each word in a sentence in every case. The tag assigned by a tagger is the most likely for a particular use of word in a sentence.

The collection of tags used by a particular tagger is called a tag set. Most part-of-speech tag sets make use of the same basic categories, i.e., noun, verb, adjective, and prepositions. However, tag sets differ in how they define categories and how finely they divide words into categories. For example, both *eat* and *eats* might be tagged as a verb in one tag set, but assigned distinct tags in another tag set. In addition, most tag sets capture morpho-syntactic information such as singular/plural, number, gender, tense, etc. Consider the following sentences:

Zuha eats an apple daily.

acquire new members while closed word classes do not (or only infrequently do so). Nouns, verbs (except auxiliary verbs), adjectives, adverbs, and interjections are open word classes. Prepositions, auxiliary verbs, delimiters, conjunction, and particles are closed word classes.

Aman ate an apple yesterday.
They have eaten all the apples in the basket.

I like to eat guavas.

The word *eat* has a distinct grammatical form in each of these four sentences. *Eat* is the base form, *ate* its past tense, and the form *eats* requires a third person singular subject. Similarly, *eaten* is the past participle form and cannot occur in another grammatical context. It is required after have or has. Thus, the following sentences are ungrammatical:

I like to eats guava.

They eaten all the apples.

The number of tags used by different taggers varies substantially. Some use 20, while others use over 400 tags. The Penn Treebank tag set contains 45 tags while C7 uses 164. For a language like English, which is not morphologically rich, the C7 tagset is too big. The tagging process would yield too many mistagged words and the result would have to be manually corrected. Despite this, bigger tag sets have been used, e.g., TOSCA-ICE for the International Corpus of English with 270 tags (Garside 1997), or TESS with 200 tags. The larger the tag set, the greater the information captured about a linguistic context. However, the task of tagging becomes complicated and requires manual correction. A bigger tag set can be used for morphologically rich languages without introducing too many tagging errors. A tag set that uses just one tag to denote all the verbs will assign identical tags to all the forms of a verb. Although this coarse-grained distinction may be appropriate for some tasks, a fine-grained tag set captures more information. This is useful for tasks like syntactic pattern detection. The Penn Treebank tag set captures finer distinctions by assigning distinct tags to distinct grammatical forms of a verb, as summarized in Table 3.9. Tags assigned to the four different forms of the word *eat* according to this tag set is shown in Table 3.10.

Table 3.9 Tags from Penn Treebank tag set

VB	Verb, base form
VBD	Subsumes imperatives, infinitives, and subjunctives
Verb, past tense	
VBG	Includes the conditional form of the verb <i>to be</i>
Verb, gerund, or present participle	
VBN	Verb, past participle
VBP	Verb, non-3rd person singular present
VBZ	Verb, 3rd person singular present

Table 3.10 Possible tags for the word to eat

eat	VB
ate	VBD
eaten	VBN
eats	VBP

Here is an example of a tagged sentence:
Speech/NN sounds/VBD sampled/VBN by/IN a/DT microphone/NN.

The tag set used is Penn Treebank.
Another tagging possible for this sentence is as follows:

Speech/NN sounds/VBZ were/VBD sampled/VBN by/IN a/DT microphone/NN.

It is easy to see that the second tagged sequence is not correct. It leads to semantic incoherence. We resolve the ambiguity using the context of the word. The context is also utilized by automatic taggers.

Part-of-speech tagging is an early stage of text processing in many NLP applications including speech synthesis, machine translation, information retrieval, and information extraction. In information retrieval, part-of-speech tagging can be used for indexing (for identifying useful tokens like nouns and phrases) and for disambiguating word senses. Tagging is not as complex as parsing. In tagging, a complete parse tree is not built; part-of-speech is assigned to words using contextual information.

Part-of-speech tagging methods fall under the three general categories.

- Rule-based (linguistic)
- Stochastic (data-driven)
- Hybrid

Rule-based taggers use hand-coded rules to assign tags to words. These rules use a lexicon to obtain a list of candidate tags and then use rules to discard incorrect tags.

Stochastic taggers have data-driven approaches in which frequency-based information is automatically derived from corpus and used to tag words. Stochastic taggers disambiguate words based on the probability that a word occurs with a particular tag. The simplest scheme is to assign the most frequent tag to each word. An early example of stochastic tagger was CLAWS (constituent likelihood automatic word-tagging system). CLAWS is the Stochastic equivalent of TAGGIT. Hidden Markov model (HMM) is the standard Stochastic tagger.

Hybrid taggers combine features of both these approaches. Like rule-based systems, they use rules to specify tags. Like stochastic systems, they use machine-learning to induce rules from a tagged training corpus automatically. The transformation-based tagger or Brill tagger is an example of the hybrid approach.

3.7.1 Rule-based Tagger

Most rule-based taggers have a two-stage architecture. The first stage is simply a dictionary look-up procedure, which returns a set of potential tags (parts-of-speech) and appropriate syntactic features for each word. The second stage uses a set of hand-coded rules to discard contextually illegitimate tags to get a single part-of-speech for each word. For example, consider the noun-verb ambiguity in the following sentence:

The show must go on.

The tagger's best guess for the word *show* in this sentence is {VB NN}. We

Resolve this ambiguity by using the following rule.

This rule simply disallows verbs after a determiner. Using this rule the word *show* in the given sentence can only be noun.

In addition to contextual information, many taggers use morphological information to help in the disambiguation process. An example of a rule that make use of morphological information is:

(VB).

Capitalization information can be utilized in the tagging of unknown

Capitalization information can be utilized in the tagging of unknown nouns. Rule-based taggers usually require supervised training. Instead, rules can be induced automatically. To induce rules untagged text is run through a tagger. The output is then manually corrected. The corrected text is then submitted to the tagger, which learns correction rules by comparing the two sets of data. This process may be repeated several times.

The earlier systems for automatic tagging were all rule-based. An example is TAGGIT (Greene and Rubin 1971), which was used for the initial tagging of the Brown corpus (Francis and Kučera 1982). This was also rule-based system. It used 3,300 disambiguation rules and was able to tag 77% of the words in the Brown corpus with their correct part-of-speech. The rest was done manually over several years. Yet another rule-based tagger is ENGTWOL (Voutilainen 1995).

3.7.2 Stochastic Tagger

The standard stochastic tagger algorithm is the HMM tagger. A Markov model applies the simplifying assumption that the probability of a chain of symbols can be approximated in terms of its parts or n -grams. The simplest n -gram model is the unigram model, which assigns the most likely tag (part-of-speech) to each token.

The unigram model needs to be trained using a tagged training corpus before it can be used to tag data. The most likely statistics are gathered over the corpus and used for tagging. The context used by the unigram tagger is the text of the word itself. For example, it will assign the tag `JJ` for each occurrence of *fast*, since *fast* is used as an adjective more frequently than it is used as a noun, verb, or adverb. This results in incorrect tagging in each of the following sentences:

She had a fast

Muslims fast during Ramadan

Those who were injured in the

In the first sentence, *fast* is used as a noun. In the second, it is a verb, and in the third, an adverb.

We would expect more ~~adverbs~~, into account when making a tagging decision. A bigram tagger uses the current word and the tag of the previous word in the tagging process. As the tag sequence “DT NN” is more likely than the tag sequence “DT JJ”, a bigram model will assign a correct tag to the word *fast* in sentence (3.1). Similarly, it is more likely that an adverb (rather than a noun or an adjective) follows a verb. Hence, in sentence (3.3), the tag assigned to *fast*

In general, an n -gram model considers the current word and the tag of the previous $n-1$ words in assigning a tag to a word. The context considered by a tri-gram model is shown in Figure 3.15. The area shaded in grey represents the context.

Tokens	w_{n-2}	w_{n-1}	w_n	w_{n+1}
tags	t_{n-2}	t_{n-1}	t_n	t_{n+1}

Figure 3.15 Context used by a tri-gram tagger

So far, we have considered how a tag is assigned to a word given the previous tag(s). However, the objective of a tagger is to assign a tag sequence to a given sentence. We now discuss how the HMM tagger assigns the most likely tag sequence to a given sentence. We call this the HMM because it uses two layers of states: a visible layer corresponding to the input words, and a hidden layer learnt by the system corresponding to the tags. While tagging the input data, we only observe the words—the tags (states) are hidden. States of the model are visible in training, not during the tagging task.

As discussed earlier, the HMM makes use of lexical and bi-gram probabilities estimated over a tagged training corpus in order to compute the most likely tag sequence for each sentence. One way to store the statistical information is to build a probability matrix. The probability matrix contains both the probability that an individual word belongs to a word class as well as the n -gram analysis, e.g., for a bigram model, the probability that a word of class X follows a word of class Y. This matrix is then used to drive the HMM tagger while tagging an unknown text.

We now return to the original problem. Given a sequence of words (sentence), the objective is to find the most probable tag sequence for the sentence.

Let W be the sequence of words.

$$W = w_1, w_2, \dots, w_n$$

The task is to find the tag sequence

$$T = t_1, t_2, \dots, t_n$$

which maximizes $P(T|W)$, i.e.,

$$T' = \operatorname{argmax}_T P(T|W)$$

Applying Bayes Rule, $P(T|W)$ can be estimated using the expression:

$$P(T|W) = P(W|T) * P(T)/P(W)$$

As the probability of the word sequence, $P(W)$, remains the same for each tag sequence, we can drop it. The expression for the most likely tag sequence becomes:

$$T' = \operatorname{argmax}_T P(W|T) * P(T)$$

Using the Markov assumption, the probability of a tag sequence can be estimated as the product of the probability of its constituent n -grams, i.e.,

$$P(T) = P(t_1) * P(t_2|t_1) * P(t_3|t_1, t_2) * \dots * P(t_n|t_1, \dots, t_{n-1})$$

$P(W|T)$ is the probability of seeing a word sequence, given a tag sequence. For example, it is asking the probability of seeing ‘The egg is rotten’ given ‘DT NNP VB JJ’. We make the following two assumptions:

- The probability of a word is dependent only on its tag.

Using these assumptions, we obtain

$$P(W|T) = P(w_1|t_1) * P(w_2|t_2) * \dots * P(w_n|t_n)$$

i.e.,

$$P(W|T) \approx \prod_{i=1}^n P(w_i|t_i)$$

So,

$$P(W|T) * P(T) = \prod_{i=1}^n P(w_i|t_i) * P(t_1) * P(t_2|t_1) * P(t_3|t_1, t_2) * \dots * P(t_n|t_1, \dots, t_{n-1})$$

Approximating the tag history using only the two previous tags, the transition probability, $P(T)$, becomes

$$P(T) = P(t_1) * P(t_2|t_1) * P(t_3|t_1, t_2) * \dots * P(t_n|t_{n-2}, t_{n-1})$$

Hence, $P(T|W)$ can be estimated as

$$\begin{aligned} P(W|T) * P(T) &= \prod_{i=1}^n P(w_i|t_i) * P(t_1) * P(t_2|t_1) * \dots * P(t_n|t_{n-2}, t_{n-1}) \\ &= \prod_{i=1}^n P(w_i|t_i) * P(t_1) * P(t_2|t_1) * \dots * P(t_n|t_{n-2}, t_{n-1}) \end{aligned}$$

We estimate these probabilities from relative frequencies via Maximum Likelihood Estimation.

$$P(t_i|t_{i-2}, t_{i-1}, t_i) = \frac{c(t_{i-2}, t_{i-1}, t_i)}{c(t_{i-2}, t_{i-1})}$$

$$P(w_i|t_i) = \frac{c(w_i, t_i)}{c(t_i)}$$

where $c(t_{i-2}, t_{i-1}, t_i)$ is the number of occurrences of t_i followed by t_{i-2}, t_{i-1} .

Stochastic models have the advantage of being accurate and language independent. Most stochastic taggers have an accuracy of 96–97%. The accuracy seems to be quite high but it should be noted that this is measured as a percentage of words. An accuracy of 96% means that for a sentence containing 20 words, the error rate per sentence will be $1 - 0.96^{20} = 56\%$. This corresponds to approximately one word per sentence.

One of the drawbacks of stochastic taggers is that they require a manually tagged corpus for training. Kupiec (1992), Cutting *et al.* (1992), and others have demonstrated that the HMM tagger can be trained from unannotated text. This makes it possible to use the model for languages in which a manually tagged corpus is not available. However, a tagger trained on a hand-coded corpus performs better than one trained on an unannotated text. In order to achieve good performance a tagged corpus is required.

We now consider an example demonstrating how the probability of a particular part-of-speech sequence for a given sentence can be computed.

Example 3.3 Consider the sentence

The bird can fly.

and the tag sequence

DT NNP MD VB

Using bi-gram approximation, the probability

$$P \left(\begin{array}{cccc} \text{DT} & \text{NNP} & \text{MD} & \text{VB} \\ | & | & | & | \\ \text{The} & \text{bird} & \text{can} & \text{fly} \end{array} \right)$$

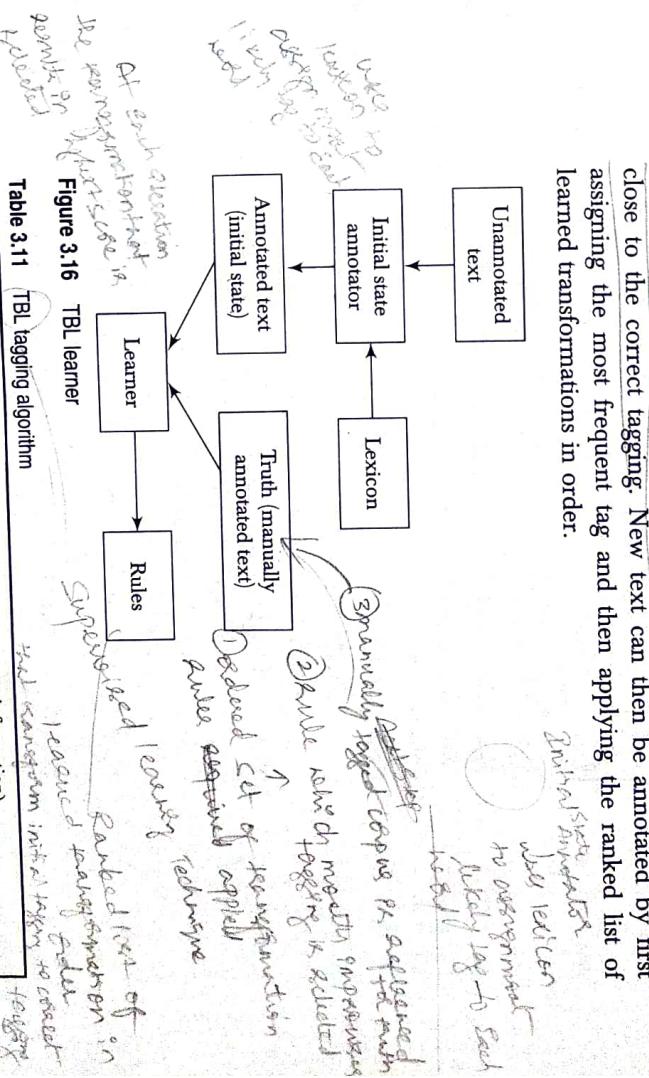
can be computed as

$$= P(\text{DT}) \times P(\text{NNP}|\text{DT}) \times P(\text{MD}|\text{NNP}) \times P(\text{VB}|\text{MD})$$

3.7.3 Hybrid Taggers

Hybrid approaches to tagging combine the features of both the rule-based and stochastic approaches. They use rules to assign tags to words like the stochastic taggers, this is a machine learning technique and rules are automatically induced from the data. Transformation-based learning (TBL) of tags, also known as Brill tagging, is an example of hybrid approach. TBL is a machine learning method introduced by E. Brill (in 1995). Transformation-based error-driven learning has been applied to a number of natural language problems, including part-of-speech tagging, speech generation, and syntactic parsing (Brill 1993, 1994, Huang *et al.* 1994).

Figure 3.16 illustrates the TBL process. Like most HMM taggers, TBL is also a supervised learning technique. The steps involved in the TBL tagging algorithm are shown in Table 3.11. The input to Brill's TBL tagging algorithm is a tagged corpus and a lexicon (with most frequent information as indicated in the training corpus). The initial state annotator uses the lexicon to assign the most likely tag to each word as the start state. An ordered set of transformation rules are applied sequentially. The rule that results in the most improved tagging is selected. A manually tagged corpus is used as reference for truth. The process is iterated until some stopping criterion is reached, such as when no significant improvement is achieved over the previous iteration. At each iteration, the transformation that results in the highest score is selected. The output of the algorithm is a ranked list of learned transformation that transform the initial tagging close to the correct tagging. New text can then be annotated by first assigning the most frequent tag and then applying the ranked list of learned transformations in order.



Each transformation is a pair of a re-write rule of the form $t_1 \rightarrow t_2$ and a contextual condition. In order to limit the set of transformations, a small set of templates is constructed. Any allowable transformation is an instantiation of these templates. Some of the transformation templates and transformations learned by TBL tagging are listed in Table 3.12.

Table 3.12 Examples of transformation templates and rules learned by the tagger

Change tag a to tag b when:

1. The preceding (following) word is tagged z.
2. The preceding (following) word is w.
3. The word two before (after) is w.
4. One of the preceding two words is w.
5. One of the two preceding (following) words is tagged z.
6. The current word is w and the preceding (following) word is x.
7. One of the previous three words is tagged z.

#	Change tags from to	Contextual condition	Example
1.	NN VB	The previous tag is TO.	To/TO fish/NN
2.	JJ RB	The previous tag is Vbz.	runs/Vbz fast/JJ

We now explain how the rules are applied in TBL tagger with the help of an example.

Example 3.4 Assume that in a corpus, *fish* is most likely to be a noun.

$$P(\text{NN}/\text{fish}) = 0.91$$

Now consider the following two sentences and their initial tags.

$$\begin{aligned} & \text{I/PRP like/VB to/TO eat/VB fish/NNP.} \\ & \text{I/PRP like/VB to/TO fish/NNP.} \end{aligned}$$

As the most likely tag for *fish* is NNP, the tagger assigns this tag to the word in both sentences. In the second case, it is a mistake.

After initial tagging when the transformation rules are applied, the tagger learns a rule that applies exactly to this mis-tagging of *fish*:

Change NNP to VB if the previous tag is TO.
As the contextual condition is satisfied, this rule will change fish/NN to fish/VB:

like/VB to/TO fish/NN \rightarrow like/VB to/TO fish/VB

The algorithm can be made more efficient by indexing the words in a training corpus using potential transformation. Recent works have involved the use of finite state transducers to compile pattern-action rules, combining

them to yield a single transducer representing the simultaneous application of all rules. Roche and Schabes (1997) have applied this approach to Brill's tagger. The resulting tagger is larger than Brill's original tagger and significantly faster.

Most of the work in part-of-speech tagging is done for English and some European languages. In other languages, part-of-speech tagging, and NLP research in general, is constrained by the lack of annotated corporuses. This is true for Indian language as well. A few part-of-speech tagging systems reported in recent years use morphological analysers along with a tagged corpus, e.g. a Bengali tagger based on HMM developed by Sandipan et al. (2004) and a Hindi tagger developed by Smriti et al. (2006). Smriti et al. used a decision tree based learning algorithm. A number of other part-of-speech taggers for Hindi, Bengali, and Telugu were developed as a result of the NLPAI-2006 machine learning contest on part-of-speech and chunking for Indian languages.

Tagging Urdu is more difficult. A number of factors contribute to this complexity. Among these is the right to left directionality of the written script and the presence of grammatical forms borrowed from Arabic and Persian. Little had been done to develop an extensive tag set for Urdu before Hardie (2003). His work was a part of the EMILLE—Enabling Minority Language Engineering project (see <http://www.emille.lancs.ac.uk/about.php>)—which focuses on the development of corpus and tools for South Asian languages.

3.7.4 Unknown Words

Unknown words are words that do not appear in dictionary or a training corpus. They create a problem during tagging. There are several potential solutions to this problem. One is to assign the most frequent tag (which occurs with most word types in the training corpus) to the unknown word. Another solution is to assume that the unknown words can be of any part-of-speech and initialize them by assigning them open class tags. Then proceed to disambiguate them using the probabilities of those tags. We can also use morphological information, such as affixes, to guess the possible tag of an unknown word. In this approach, the unknown word is assigned a tag based on the probability of the words belonging to a specific part-of-speech in the training corpus having the same suffix or prefix. A similar approach is used in Brill's tagger.

MARY

- This chapter has dealt with word level analysis. The topics covered include methods for characterizing word sequences, identifying morphological variants, detecting and correcting misspelled words, and identifying the correct part-of-speech for a word. The main points are as follows.
- Regular expressions can be used for specifying words. They can be encoded as a finite automaton.
- The goal of morphological parsing is to find out morphemes using which a given word is built. Morphemes are the smallest meaning bearing units in a language.
- Morphological analysis and generation are essential to many NLP applications ranging from spelling error corrections to machine translations.
- The simplest morphological systems are stemmers. They do not use a lexicon. Instead, they use re-write rules. However, stemmers are not perfect.
- A two-level morphological model is more efficient. Both of its steps can be implemented using a finite state transducer.
- Word errors belong to one of two distinct categories, namely, *non-word errors* and *real word errors*. The latter category requires the context of the word to detect and correct errors.
- Words are classified into categories called part-of-speech or word classes. Word classes can be open or closed.
- Part-of-speech tagging is the process of assigning a part-of-speech like noun, verb, pronoun, preposition, adverb, adjective, etc., to each word in a sentence.
- Part-of-speech tagging methods fall under the following three categories:
 1. Rule-based (linguistic)
 2. Stochastic (data-driven)
 3. Hybrid
- *Rule-based taggers* use hand-coded rules to assign tags to words. *Stochastic taggers* require a pre-tagged corpus for training. *Hybrid taggers* combine features of both these approaches. Like rule-based systems, they use rules to specify tags. Like stochastic methods, they use machine learning to automatically induce rules from a tagged training corpus.
- Unknown words can be assigned the most frequent tags. We can also use morphological information to guess the correct tag.

REFERENCES

- Brill, E., 1993, 'Transformation-based error-driven parsing,' *Proceedings of the Third International Workshop on Parsing Technologies*, Tilburg, The Netherlands.
- Brill, E., 1994, 'Some advances in rule-based part of speech tagging,' *Proceedings of the Twelfth National Conference on Artificial Intelligence (AAAI-94)*, Seattle.
- Brill, E., 1995, 'Error-driven learning and natural language processing: a case study in part-of-speech tagging,' *Computational Linguistics*.
- Cutting, D., J. Kupiec, J. Pederson, and P. Sibun, 1992, 'A practical part-of-speech tagger,' *Proceedings of the Third Conference on Applied Natural Language Processing, ACL*.
- Damerau, F.J., 1964, 'A technique for computer detection and correction of spelling errors,' *Communications of the ACM*, 7(3) pp. 171–76.
- Francis, W. Nelson and Henry Kucera, 1982, *Frequency Analysis of English Usage: Lexicon and Grammar*, Houghton Mifflin, Boston.
- Garside, R., G. Leech, and G. Sampson, 1987, *The Computational Analysis of English: A Corpus-Based Approach*, Longman, London.
- Green, B. and G. Rubin, 1971, 'Automated grammatical tagging of English,' Department of Linguistics, Brown University.
- Hardie, A., 2003, 'Developing a tag-set for automated part-of-speech tagging in Urdu,' *Proceedings of the Corpus Linguistics 2003 Conference*.
- D. Archer, P. Rayson, A. Wilson, and T. McEnery, (Eds.), UCREL Technical Papers, 16, Department of Linguistics, Lancaster University.
- Hopcroft, J.E. and J.D. Ullman, 1979, *Introduction to Automata Theory, Languages, and Computation*, Addison-Wesley, Reading, Massachusetts.
- Huang et al., 1994, 'Generation of pronunciation from orthographies using transformation-based error-driven learning,' *Proceedings of Int. Conference on Speech and Language Processing (ICSLP)*, Yokohama, Japan.
- Kleene, S.C., 1956, 'Representation of events in nerve nets and finite automata,' *Automata Studies*, C. Shannon and J. McCarthy (Eds.), Princeton University Press, Princeton, NJ, pp. 3–41.
- Koskeniemi, K., 1983, 'Two-level morphology: A general computational model of word-form recognition and production, Technical Report Publication No. 17, Department of General Linguistics,' University of Helsinki.
- Krovetz, R., 1993, 'Viewing morphology as an inference process,' In SIGIR-93, pp. 191–202.

- Kukich, K., 1992, 'Techniques for automatically correcting words in text,' *ACM Computing Surveys*, 24, pp. 377–439.
- Kupiec, J., 1992, 'Robust part-of-speech tagging using a Hidden Markov Model,' *Computer Speech and Language*, vol. 6, pp. 225–42.
- Lovins, J. B., 1968, 'Development of a stemming algorithm,' *Mechanical Translation and Computational Linguistics*, 11(1–2), pp. 22–31.
- Odell, M. K. and R. C. Russell, US Patents 1261167/1435663 (1918/1922).
- Oflazer, Kemal, 1996, 'Error-tolerant finite state recognition with applications to morphological analysis and spelling correction,' *Computational Linguistics*, 22(1).
- Porter, M. F., 1980, 'An algorithm for suffix stripping program,' 14(3), pp. 130–37.
- Roche, E. and Y. Schabes, 1995, 'Deterministic part of speech tagging with finite state transducers,' *Computational Linguistics*.
- Sandipan, D., Kumar Nagraj, and Uma Sawant, 2004, 'A hybrid model for part-of-speech tagging and its application to Bengali,' *Proceedings of International Conference on Computational Intelligence*.
- Shaffer, L. and J. Hardwick, 1968, 'Typing performance as a function of text,' *Quarterly Journal of Experimental Psychology*, 20, pp. 360–69.
- Singh, Smriti, Kuhoo Gupta, Manish Shrivastava, and Pushpak Bhattacharyya, 2006, 'Morphological richness offsets resource demand-experiences in constructing a POS tagger for Hindi,' *Proceedings of the COLING/ACL 2006 Main Conference Poster Sessions*, Association for Computational Linguistics, Sydney, pp. 779–86.
- Voutilainen, A., 1996, 'Morphological disambiguation,' *Constraint Grammar: A language-independent system for parsing uncertainty text*, F. Karlsson, A. Voutilainen, J. Heikkila, and A. Anttila (Eds.), Berlin, pp. 165–284.
- Wagner, R.A. and M.J. Fischer, 1974, 'The string-to-string correction problem,' *Journal of the Association for Computing Machinery*, 21, 168–73.

EXERCISES

- Define a finite automaton that accepts the following language:
(aa)*(bb)*.
- A typical URL is of the form:

http	://	www.abc.com	/mfpaper/public	/xxx.html
1	2	3	4	5

LAB EXERCISES

- Write a program to find minimum edit distance between two input strings.
- Use any tagger available in your lab to tag a text file. Now write a program to find the most likely tag in the tagged text.
- Write a program to find the probability of a tag given previous two tags, i.e., $P(t_3/t_2 t_1)$.
- How can unknown words be handled in the tagging process?

In this table, 1 is a protocol, 2 is name of a server, 3 is the directory, and 4 is the name of a document. Suppose you have to write a program that takes a URL and returns the protocol used, the DNS name of the server, the directory and the document name. Develop a regular expression that will help you in writing this program.

- Distinguish between non-word and real-word error.
- Compute the minimum edit distance between *paeffu* and *peaceful*.
- Comment on the validity of the following statements:
 - Rule-based taggers are non-deterministic.
 - Stochastic taggers are language independent.
 - Brill's tagger is a rule-based tagger.

SYNTACTIC ANALYSIS

Researchers have proposed a number of parsing methods for natural language sentences. It is not possible to discuss all of them in a chapter, so we focus on a few widely-known ones. Most text books use only English for their discussions. We differ by including Indian languages as well. In addition to an introduction to grammar formalism, this chapter also provides an introduction to Hindi grammar.

Two important ideas in natural language are those of constituency and word order. Constituency is about how words are grouped together and how we know that they are really grouping together. Word order is about how, within a constituent, words are ordered with respect to one another, and also how constituents are ordered with respect to one another.

A widely used mathematical system for modelling constituent structure in natural language is context-free grammar (CFG) also known as phrase structure grammar. We begin our discussion with an overview of CFG in Section 4.2. In Section 4.3, we discuss various types of phrases and phrase structure rules. We then introduce feature structures in Section 4.4, to capture certain properties of grammatical categories which cannot be handled efficiently using CFG. We discuss probabilistic grammar and CYK parser based on this in Section 4.5. The CFG is basically a positional grammar and is not suitable for Indian languages, which are free word order languages. We provide a brief overview of Paninian grammar (PG), which is suitable for modelling free word order languages, in Section 4.6. We also discuss a parsing framework proposed by Bharti and Sangal (1990) for Indian languages based on this grammar. Finally, we provide a brief summary of the chapter.

4.1 INTRODUCTION

In Chapter 3, we talked about word level analysis. We now move on to higher level constituents like phrases and sentences. The word ‘syntax’ refers to the grammatical arrangement of words in a sentence and their relationship with each other. The objective of syntactic analysis is to find the syntactic structure of the sentence. This structure is usually depicted as a tree, as shown in Figure 4.1. Nodes in the tree represent the phrases and leaves correspond to the words. The root of the tree is the whole sentence. Identifying the syntactic structure is useful in determining the meaning of the sentence. The identification is done using a process known as parsing. Syntactic parsing can be considered as the process of assigning ‘phrase markers’ to a sentence (Charniak 1997). One must, therefore, have a clear understanding of what phrases are. We will describe various phrases that could appear in a language so that we can specify grammar rules for them.

4.2 CONTEXT-FREE GRAMMAR

Context-free grammar (CFG) was first defined for natural language by Chomsky (1957) and used for the Algol programming language by Backus (1959) and Naur (1960). A CFG (also called phrase-structure grammar) consists of four components:

1. A set of non-terminal symbols, N
2. A set of terminal symbols, T
3. A designated start symbol, S , that is one of the symbols from N .
4. A set of productions, P , of the form:

$$A \rightarrow \alpha$$

where $A \in N$ and α is a string consisting of terminal and non-terminal symbols. The rule $A \rightarrow \alpha$ says that constituent A can be rewritten as α . This is also called the phrase structure rule. It specifies which elements

4.3 CONSTITUENCY

(or constituents) can occur in a phrase and in what order. For example, the rule $S \rightarrow NP\ VP$ states that S consists of NP followed by VP , i.e., a sentence consists of a noun phrase followed by a verb phrase.

A language is usually defined through the concept of derivation. The basic operation is that of rewriting a symbol appearing on the left hand side of production by its right hand side. A CFG can be used to generate a sentence or to assign a structure to a given sentence. When used as a generator, the arrows in the production rule may be read as ‘rewrite the symbol on the left with symbols on the right’. Consider the toy grammar shown in Figure 4.1. The symbol S can be rewritten as $NP\ VP$ using Rule 1, then using rules R2 and R4, NP and VP are rewritten as N and V respectively. NP is then rewritten as $Det\ N$ (R3). Finally, using rules R6 and R7, we get the sentence:

Hena reads a book.

We say that the sentence (4.1) can be derived from S . The representation of this derivation is shown in Figure 4.1. Sometimes, a more compact bracketed notation is used to represent a parse tree. The parse tree in Figure 4.1 can be represented using this notation as follows:

$[S [NP [N Hena]] [VP [V reads] [NP [Det a] [N book]]]]]$

The set of all the strings containing terminal symbols which can be derived from the start symbol of the grammar, defines the language generated by that grammar. The parse tree shown in Figure 4.1 essentially represents a mapping of a string to its parse tree. This mapping process is called parsing. We will come back to this issue in Section 4.4.

Words in a sentence are not tied together as a sequence of part-of-speech. Language puts constraints on word order. For example, certain words go together with each other more than with others, and seem to behave as a unit. The fundamental idea of syntax is that words group together to form constituents (often termed phrases), each of which acts as a single unit. They combine with other constituents to form larger constituents, and eventually, a sentence. *The bird*, *The rain*, *The Wimbledon court*, *The beautiful garden* are all noun phrases that can occur in the same syntactic context. For example, they can all function as the subject or the object of a verb. These constituents combine with others to form a sentence constituent. For example, the noun phrase, *The bird*, can combine with the verb phrase, *flies*, to form the sentence, *The bird flies*. Different types of phrases have different internal structures. In this section, we discuss some of the major phrase types and try to build phrase structure rules to identify them.

4.3.1 Phrase Level Constructors

As discussed earlier, a fundamental notion in natural language is that certain groups of words behave as constituents. These constituents are identified by their ability to occur in similar contexts. One of the simplest ways to decide whether a group of words is a phrase, is to see if it can be substituted with some other group of words without changing the meaning. If such a substitution is possible then the set of words forms a phrase. This is called the substitution test. Consider sentence (4.1). We can substitute a number of other phrases:

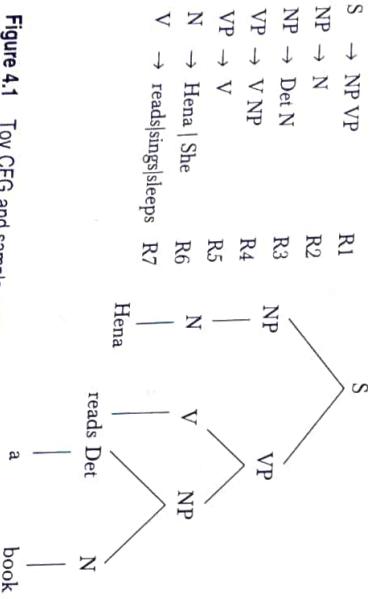
Hena reads a book.
Hena reads a storybook.
Those girls read a book.

She reads a comic book.

We can easily identify the constituents that can be replaced for each other in these sentences. These are *Hena*, *she*, and *Those girls* and *a book*, *a storybook*, and *a comic book*. These are the words that form a phrase. In linguistics, such constituents represent a paradigmatic relationship. Elements that can substitute each other in certain syntactic positions are said to be members of one paradigm.

Phrase types are named after their head, which is the lexical category

Figure 4.1 Toy CFG and sample parse tree



called a verb phrase; and so on for other lexical categories such as adjective and preposition. Figure 4.2 shows a sentence with a noun phrase, verb phrase, and preposition phrase.

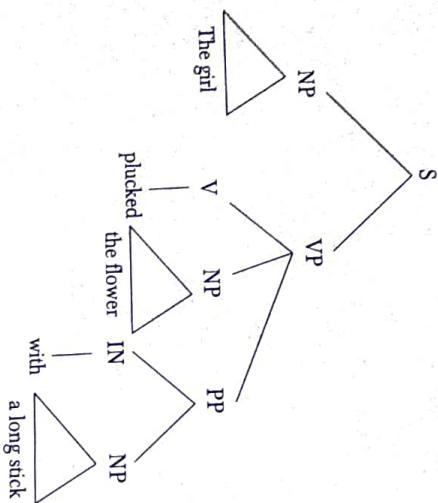


Figure 4.2 A sentence with NP, VP, and PP

Noun Phrase

A noun phrase is a phrase whose head is a noun or a pronoun, optionally accompanied by a set of modifiers. It can function as subject, object, or complement. The modifiers of a noun phrase can be determiners or adjective phrases. The obligatory constituent of a noun phrase is the noun head—all other constituents are optional. These structures can be represented using the phrase structure rule. As discussed earlier, phrase structure rules are of the form $A \rightarrow BC$, which states that constituent A can be rewritten as two constituents B and C. These rules specify which elements can occur in a phrase and in what order. Using this notation, we can represent the phrase structure rules for a noun phrase as follows.

$$\begin{aligned} \text{NP} &\rightarrow \text{Pronoun} \\ \text{NP} &\rightarrow \text{Det Noun} \\ \text{NP} &\rightarrow \text{Noun} \\ \text{NP} &\rightarrow \text{Adj Noun} \\ \text{NP} &\rightarrow \text{Det Adj Noun} \end{aligned}$$

We can combine all these rules in a single phrase structure rule as follows:

$$\text{NP} \rightarrow (\text{Det}) (\text{Adj}) \text{Noun} | \text{Pronoun}$$

The constituents in parentheses are optional. This rule states that a noun phrase consists of a noun, possibly preceded by a determiner and an adjective (in that order). This rule does not cover all possible NPs. A

The following are a few examples of noun phrases:

The foggy morning (4.2a)

Chilled water (4.2b)

A beautiful lake in Kashmir (4.2c)

Cold banana shake (4.2d)

Their new house (4.2e)

Let us see how the phrases (4.2a–e) can be generated using phrase structure rules. The phrase (4.2a) consists only of a noun; (4.2b) consists of a determiner, an adjective (foggy) that stands for an entire adjective phrase, and a noun; (4.2c) comprises an adjective phrase and a noun; (4.2d) consists of a determiner (the), an adjective phrase (beautiful), a noun (lake), and a prepositional phrase (in Kashmir); and (4.2e) consists of an adjective followed by a sequence of nouns. A noun sequence is termed as nominal. None of the phrase structure rules discussed so far are able to handle nominals. So, we modify our rules to cover this situation.

$$\begin{aligned} \text{NP} &\rightarrow (\text{Det}) (\text{AP}) \text{ Nom} | (\text{PP}) \\ \text{Nom} &\rightarrow \text{Noun} | \text{Noun Nom} \end{aligned}$$

A noun phrase can act as a subject, an object, or a predicate. The following sentences demonstrate each of these uses.

The foggy damped weather disturbed the match. (4.3a)

I would like a nice cold banana shake. (4.3b)

Kula botanical garden is a beautiful location. (4.3c)

In (4.3a), the noun phrase acts as a subject. In (4.3b), it acts as an object, and in (4.3c), it is a predicate.

Verb Phrase

Analogous to the noun phrase is the verb phrase, which is headed by a verb. There is a fairly wide range of phrases that can modify a verb. This makes verb phrases a bit more complex. The verb phrase organizes various elements of the sentence that depend syntactically on the verb.

The following are some examples of verb phrases:

$$\begin{aligned} (4.4a) \quad &\text{The boy kicked the ball.} \\ (4.4b) \quad &\text{Khushbu slept.} \end{aligned}$$

Khushbu slept in the garden.

The boy gave the girl a book.

The boy gave the girl a book with blue cover.

(4.4c)
(4.4d)

(4.4e)

4.3.2 Sentence Level Constructions

As you can see from these examples a verb phrase can have a verb [VP → Verb in (4.4a)]; a verb followed by an NP [VP → Verb NP in (4.4b)]; a verb followed by a PP [VP → Verb PP in (4.4c)]; a verb followed by two NPs [VP → Verb NP NP in (4.4d)]; or a verb followed by two NPs and a PP [VP → Verb NP NP PP in (4.4e)]. In general, the number of NPs in a VP is limited to two, whereas it is possible to add more than two PPs.

$\text{VP} \rightarrow \text{Verb (NP) (NP) (PP)}^*$

Things are further complicated by the fact that objects may also be entire clauses as in the sentence, *I know that Taj is one of the seven wonders*. Hence, we must also allow for an alternative phrase statement rule, in which NP is replaced by S.

$\text{VP} \rightarrow \text{Verb S}$

Prepositional Phrase

Prepositional phrases are headed by a preposition. They consist of a preposition, possibly followed by some other constituent, usually a noun phrase.

We played volleyball on the beach.

We can have a preposition phrase that consists of just a preposition.

John went outside.

The phrase structure rule that captures the above eventualities is as follows.

$\text{PP} \rightarrow \text{Prep (NP)}$

The head of an adjective phrase (AP) is an adjective. APs consist of an adjective, which may be preceded by an adverb and followed by a PP.

Adjective Phrase

Here are few examples.

Ashish is clever.

The train is very late.

My sister is fond of animals.

The phrase structure rule for adjective phrase is

$\text{AP} \rightarrow (\text{Adv}) \text{ Adj (PP)}$

Adverb Phrase

An adverb phrase consists of an adverb, possibly preceded by a degree adverb. Here is an example.

Time passes very quickly.
 $\text{AdvP} \rightarrow (\text{Intens}) \text{ Adv}$

Sentences with an imperative structure usually begin with a verb phrase and lack subject. The subject of these types of sentence is implicit and is understood to be ‘you’. These types of sentences are used for commands and suggestions, and hence are called imperative. The grammar rule for this kind of sentence structure is

$S \rightarrow \text{NP VP}$

Sentences with the yes-no question structure ask questions which can be answered using yes or no. These sentences begin with an auxiliary verb, followed by a subject NP, followed by a VP. Here are some examples:

Do you have a red pen?

Is there a vacant quarter?

Is the game over?

Can you show me your album?

We expand our grammar by adding another rule for the expansion of S, as follows:

$S \rightarrow \text{Aux NP VP}$

Sentences with wh-question structure are more complex. These sentences begin with a wh-words—who, which, where, what, why, and how. A wh-question may have a wh-phrase as a subject or may include another subject. Consider the following wh-question:

Which team won the match?

This sentence is similar to a declarative sentence except that it contains a wh-word. A simple rule to handle this type of sentence structure is

$$S \rightarrow Wh\text{-NP VP}$$

Another type of wh-question structure is one that involves more than one NP. In this type of questions, the auxiliary verb comes before the subject NP, just as in yes-no question structures.

Which cameras can you show me in your shop?

The rule for this type of wh-questions is

$$S \rightarrow Wh\text{-NP Aux NP VP}$$

A simplified view of the grammar rules discussed so far is summarized in Table 4.1.

Table 4.1 Summary of grammar rules

Dr. R.R.C.K.R.	DBIT
Acc. No. 44954	

Coordination

The grammar rules in Table 4.1 are not exhaustive. There are other sentence-level structures that cannot be modelled by the rules discussed here. Coordination is one such structure. It refers to conjoining phrases with conjunctions like ‘and’, ‘or’, and ‘but’. For example, a coordinate noun phrase can consist of two other noun phrases separated by a conjunction ‘and’, as in

I ate [NP [NP an apple] and [NP a banana]].

Similarly, verb phrases and prepositional phrases can be conjoined as follows:

It is [VP [VP dazzling] and [VP raining]].

Not only that, even a sentence can be conjoined.

[S [S I am reading the book] and [S I am also watching the movie]]

We need to devise rules to handle these constructions. Conjunction rules for NP, VP, and S can be built as follows:

$$NP \rightarrow NP \text{ and } NP$$

$$VP \rightarrow VP \text{ and } VP$$

$$S \rightarrow S \text{ and } S$$

$SgNP$	$\rightarrow (Det) (AP) SgNom (PP)$
$NonSgNP$	$\rightarrow (Det) (AP) PlNom (PP)$
$SgNom$	$\rightarrow SgNom SgNom SgNom$
$PlNom$	$\rightarrow PlNom PlNom PlNom$
$SgNoun$	$\rightarrow Priya lake banana sister \dots$
$PlNoun$	$\rightarrow Children \dots$

We also have to add rules for the first and second person pronouns. This method of dealing with number agreement doubles the size of the grammar. Each rule that makes use of a noun or verb phrase results in the introduction of a pair of rules—one to handle singular form and another to handle plural form. We also need to introduce new versions of NP and noun rules for various cases, e.g., nominative (I, she, they, he) and accusative (me, her, him, them) cases of pronoun. Languages like Hindi and Urdu, which have not only noun-verb agreements but also gender agreements, further aggravate the problem by adding another multiplier. Clearly, CFG cannot handle this problem efficiently.

We solve the problem of over-generation by introducing new grammatical categories corresponding to each such constraint. This results in an explosion in the number of grammar rules and loss of generality.

An alternative solution is to associate each non-terminal of the grammar with feature structures. Feature structures are able to capture grammatical properties without increasing the size of the grammar. We can think of grammatical categories (and the grammar rules that make use of them) as objects having properties associated with them. The information in these properties can be thought of as constraints imposed by the grammatical categories. Models based on this idea are called constraint-based formalisms. We now introduce feature structures and discuss how they are used to represent the constraints imposed by grammatical categories without the loss of generality.

Feature Structures

Feature structures are sets of feature-value pairs. They can be used to efficiently capture the properties of grammatical categories. Features are simply symbols representing properties that we wish to capture. For example, the number property of a noun phrase can be represented by NUMBER feature. The value that a NUMBER feature can take is SG (for singular) and PL (for plural). Values can be either atomic symbols or feature structures. Feature structures are represented by a matrix-like diagram called attribute value matrix (AVM).

$[FEATURE_1$	$VALUE_1]$
$FEATURE_2$	$VALUE_2]$
\dots	$\dots]$
$FEATURE_n$	$VALUE_n]$

Figure 4.3 An attribute value matrix (AVM)

An AVM consisting of a single NUMBER feature with the value SG is represented as follows:

$[NUMBER SG]$

The value of a feature can be left unspecified and represented by an empty pair of square brackets, as in the following example:

$[NUMBER []]$

The feature structure can be used to encode the grammatical category of a constituent and the features associated with it. For example, the following structure represents the third person singular noun phrase.

CAT	NP
$NUMBER$	SG
$PERSON$	3

Similarly, a third person plural noun phrase can be represented as follows:

CAT	NP
$NUMBER$	PL
$PERSON$	3

The values of CAT and PERSON features remain the same in both structures. This explains how feature structures aid in generalization while making the necessary distinction possible. As mentioned earlier, feature values are not limited to atomic symbols. A feature can have another feature structure as its value. Consider the case of combining the NUMBER and PERSON features into a single AGREEMENT feature. This makes sense because grammatical subjects must agree with their predicates in NUMBER as well as PERSON properties. Using this new feature, we represent the grammatical category third person plural noun phrase by the following structure:

CAT	NP
$AGREEMENT$	$[NUMBER PL]$
$PERSON$	3

In order for feature structures to be useful, we must be able to perform operations on them. The two most important operations we need to perform are merging the information content of the two structures that are similar and rejecting structures that are incompatible. The computational technique that is used to perform these operations is called unification. Unification is implemented as a binary operator (\sqcup) that takes two feature structures as arguments and returns a merged feature structure if they are compatible, otherwise reports a failure. Here is a simple application of the unification operator for performing an equality check.

$$[\text{NUMBER } \textit{PL}] \sqcup [\text{NUMBER } \textit{PL}] = [\text{NUMBER } \textit{PL}]$$

The unification succeeds as the two structures have the same value for the NUMBER feature. A feature with an unspecified value in one structure, can be successfully matched with any value in a corresponding feature in another structure. In such cases, the unification operation produces a structure with the value provided by the structure having non-null value. For example,

$$[\text{NUMBER } \textit{PL}] \sqcup [\text{NUMBER } []] = [\text{NUMBER } \textit{PL}]$$

In this example, the two structures are considered compatible and merged into a structure with PL as its value for the NUMBER feature. The value PL of the first structure matches the value [] of the second structure and becomes the value of the NUMBER feature of the output structure.

However, the following application of unification results in failure as the NUMBER features of the first and second structures have incompatible values.

$$[\text{NUMBER } \textit{PL}] \sqcup [\text{NUMBER } \textit{SG}] \text{ Fails}$$

The CFG rules can have feature structures attached to them to realize constraints on the constituents of the sentence.

4.4 PARSING

A CFG defines the syntax of a language but does not specify how structures are assigned. The task that uses the rewrite rules of a grammar to either generate a particular sequence of words or reconstruct its derivation (or phrase structure tree) is termed parsing. A phrase structure tree constructed from a sentence is called a parse. The syntactic parser is thus responsible for recognizing a sentence and assigning a syntactic structure to it. It is possible for many different phrase structure trees to derive the same sequence of words. This means a sentence can have multiple parses. This phenomenon is called syntactic ambiguity.

Garden pathing is another phenomenon related to syntactic parsing. It refers to the process of constructing a parse by exploring the parse tree along different paths, one after the other till, eventually, the right one is found. The popular example given for garden pathing is the sentence

The horse ran past the barn fell.

In the first attempt, most of us come up with a parse corresponding to the sentence *The horse ran past the barn*, leaving no possibility for the word *fell* to be incrementally added in the sentence. In order to complete the parse of the sentence, we have to backtrack.

Finding the right parse can be viewed as a search process. The search finds all trees whose root is the start symbol S and whose leaves cover exactly the word in the input. The search space in this conception corresponds to all possible parse trees defined by the grammar. The following constraints guide the search process.

1. *Input:* The first constraint comes from the words in the input sentence. A valid parse is one that covers all the words in a sentence. Hence, these words must constitute the leaves of the final parse tree.

2. *Grammar:* The second kind of constraint comes from the grammar. The root of the final parse tree must be the start symbol of the grammar.

These two constraints give rise to the two most widely used search strategies by parsers, namely, top-down or goal-directed search and bottom-up or data-directed search.

4.4.1 Top-down Parsing

As the name suggests, top-down parsing starts its search from the root node S and works downwards towards the leaves. The underlying assumption here is that the input can be derived from the designated start symbol, S, of the grammar. The next step is to find all sub-trees which can start with S. To generate the sub-trees of the second-level search, we expand the root node using all the grammar rules with S on their left hand side. Likewise, each non-terminal symbol in the resulting sub-trees is expanded next using the grammar rules having a matching non-terminal symbol on their left hand side. The right hand side of the grammar rules provide the nodes to be generated, which are then expanded recursively. As the expansion continues, the tree grows downward and eventually reaches a state where the bottom of the tree consist only of part-of-speech categories. At this point, all trees whose leaves do not match words in the input sentence are rejected, leaving only trees that represent successful

parses. A successful parse corresponds to a tree which matches exactly with the words in the input sentence.

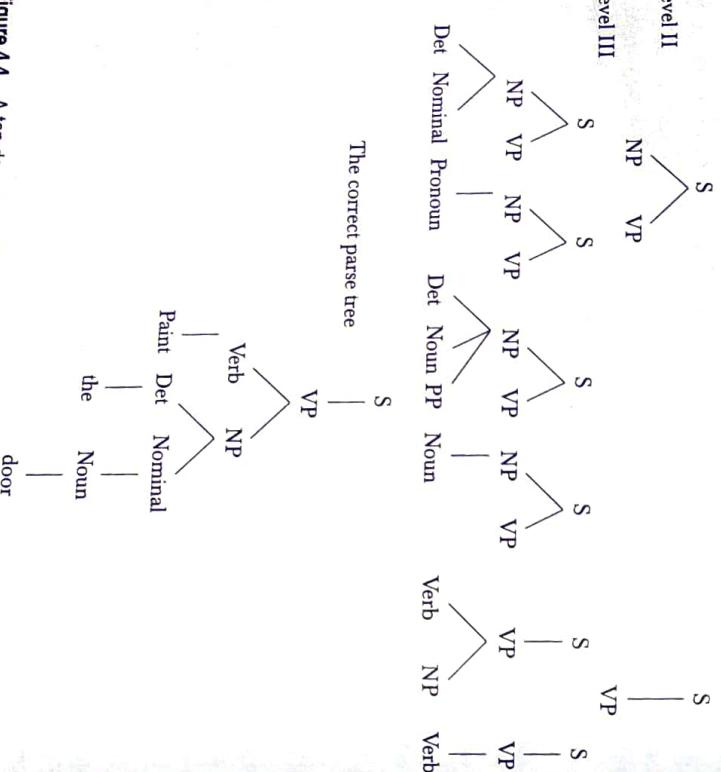
Table 4.2 Sample grammar

$S \rightarrow NP\ VP$	VP \rightarrow Verb NP
$S \rightarrow VP$	VP \rightarrow Verb
$NP \rightarrow Det\ Nominal$	PP \rightarrow Preposition NP
$Det \rightarrow this that a the$	Verb \rightarrow sleeps sings open saw paint
$NP \rightarrow Noun$	Preposition \rightarrow from with on to
$Nominal \rightarrow Noun\ Nominal$	Pronoun \rightarrow she he they
$Nominal \rightarrow Noun\ Nominal$	

Consider the grammar shown in Table 4.2 and the sentence

Paint the door.

A top-down search begins with the start symbol of the grammar. Thus, the first level (ply) search tree consists of a single node labelled S. The grammar in Table 4.2 has two rules with S on their left hand side. These



(4.7)

A bottom-up parser starts with the words in the input sentence and attempts to construct a parse tree in an upward direction towards the root. At each step, the parser looks for rules in the grammar where the right hand side matches some of the portions in the parse tree constructed so far, and reduces it using the left hand side of the production. The parse is considered successful if the parser reduces the tree to the start symbol of the grammar. Figure 4.5 shows some steps carried out by the bottom-up parser for sentence (4.7).

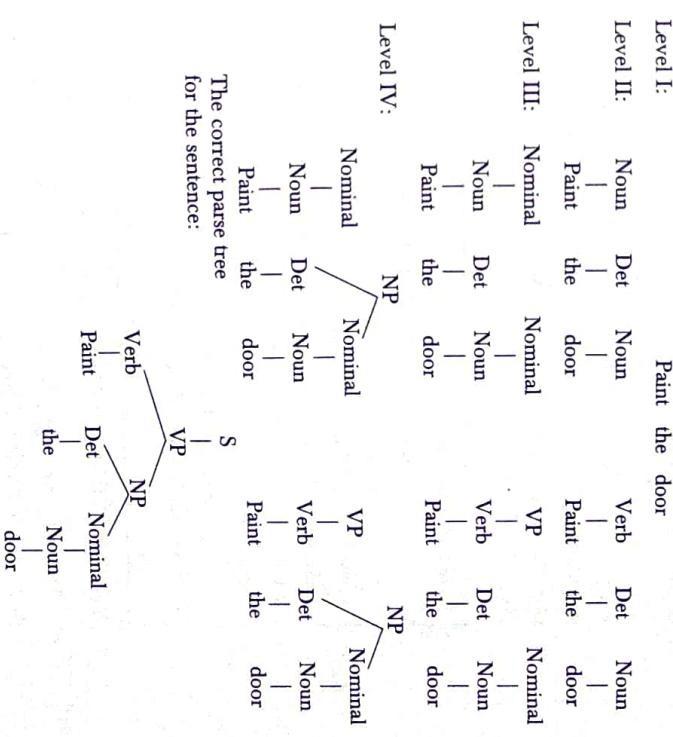


Figure 4.4 A top-down search space

Figure 4.5 A bottom-up search space for sentence (4.7)

rules are used to expand the tree, which gives us two partial trees at the second level search, as shown in Figure 4.4. The third level is generated by expanding the non-terminal at the bottom of the search tree in the previous ply. Due to space constraints, only the expansion corresponding to the left-most non-terminals has been shown in the figure. The subsequent steps in the parse are left, as an exercise, to the readers. The correct parse tree shown in Figure 4.4 is obtained by expanding the fifth parse tree of the third level.

4.4.2 Bottom-up Parsing

A bottom-up parser starts with the words in the input sentence and attempts to construct a parse tree in an upward direction towards the root. At each step, the parser looks for rules in the grammar where the right hand side matches some of the portions in the parse tree constructed so far, and reduces it using the left hand side of the production. The parse is considered successful if the parser reduces the tree to the start symbol of the grammar. Figure 4.5 shows some steps carried out by the bottom-up parser for sentence (4.7).

Each of these parsing strategies has its advantages and disadvantages.

As the top-down search starts generating trees with the start symbol of the grammar, it never wastes time exploring S trees that eventually result in words that are inconsistent with the input. This is because a top-down parser generates trees before seeing the input. On the other hand, a bottom-up parser never explores a tree that does not match the input. However, it wastes time generating trees that have no chance of leading to an S -rooted tree. The left branch of the search space in Figure 4.5 that explores a sub-tree assuming *paint* as a noun, is an example of wasted effort. We now present a basic search strategy that uses the top-down method to generate trees and augments it with bottom-up constraints to filter bad parses.

4.4.3 A Basic Top-down Parser

The approach presented here is essentially a depth first, left to right search. The depth first approach expands the search space incrementally by one state at a time. At each step, the left-most unexpanded leaf nodes of the tree are expanded first using the relevant rule of the grammar. The left-most node is selected for expansion as it determines the order in which input words needs to be considered. When a state arrives that is inconsistent with the input, the search continues by returning to the most recently generated and unexplored tree. The steps of the algorithm are given in Figure 4.6.

1. Initialize agenda
2. Pick a state, let it be curr_state, from agenda
3. If (curr_state) represents a successful parse then return parse tree
 - else if curr_state is a POS then
 - if category of curr_state is a subset of POS associated with curr_word then apply lexical rules to current state
 - else reject
 - else generate new states by applying grammar rules and push them into agenda
 4. If (agenda is empty) then return failure
 - else select a node from agenda for expansion and go to step 3.

Figure 4.6 Top-down, depth-first parsing algorithm

The algorithm maintains an agenda of search states. Each search state consists of partial trees and a pointer to the next input word in the sentence. The algorithm starts with the state at the front of the agenda

and generates a set of new states by applying grammar rule to the left-most unexpanded node of the tree associated with it. The newly generated states are put on the front of the agenda in the order defined by the textual order of the grammar rules used to create them. The process continues until either a successful parse tree is discovered or the agenda is empty, indicating a failure.

Figure 4.7 shows the trace of the algorithm on the sentence, *Open the door*. The algorithm starts with the node S and input word *Open*. It first expands S using the grammar rule $S \rightarrow NP VP$. It then expands the left-most unexpanded non-terminal NP using the rule $NP \rightarrow Det Nominal$. But the word *Open* cannot be derived from *Det*. Hence, the parser eliminates the rule and tries the second alternative, i.e., $NP \rightarrow noun$, which again leads to a failure. The next search space on the agenda corresponds to the $S \rightarrow VP$ rule. The expansion of VP using the rule $VP \rightarrow Verb NP$, successfully matches the first input words. The algorithm proceeds in a depth-first, left-to-right manner, to match the rest of the input words.

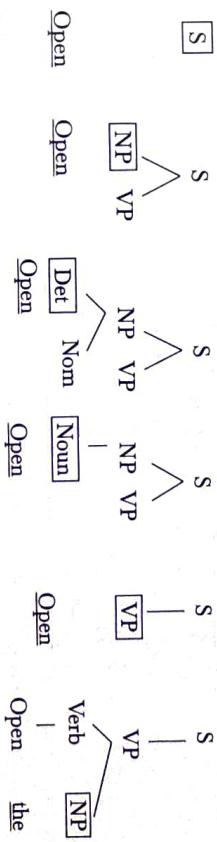


Figure 4.7 A derivation using top-down, depth-first algorithm

In any successful parse, the current input word must match the first word in the derivation of the node that is being expanded. This information can be utilized in eliminating spurious parses. A grammar rule that cannot lead to the input word as the first word along the left side of a derivation, should not be considered for expansion. The first word along the left side

of the derivation is called the left corner of the tree. Using the left corner notion, we see that in our example, only the rule $S \rightarrow VP$ is applicable, as the word *Open* cannot be a left corner of the NP. In order to utilize this filter, we create a table containing a list of all the valid left corner categories for each non-terminal of the grammar. While selecting a rule for expansion, the table is consulted to see if the non-terminal associated with the rule has a part-of-speech associated with the current input. If not, then the rule is not considered. The left corner table for grammar is shown in Table 4.3.

The top-down, depth-first, left-to-right search algorithm suffers from certain disadvantages. The first is that of left recursion, which causes the search to get stuck in an infinite loop. This problem arises if the grammar is left recursive that, is, it contains a non-terminal A, which derives, in one or more steps, a string beginning with the same non-terminal, i.e., $A^* \Rightarrow A\beta$ for some β .

Table 4.3 Left corner for each grammar category

Category	Left Corners
S	Det, Pronoun, Noun, Verb
NP	Noun, Pronoun, Det
VP	Verb
PP	Preposition
Nominal	Noun

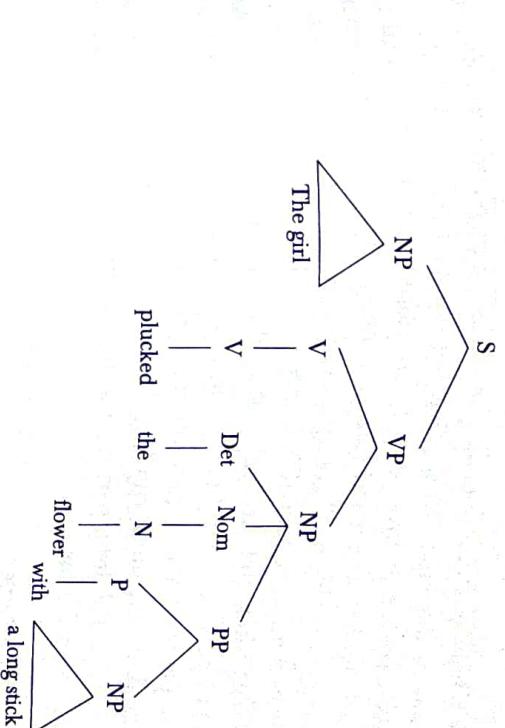


Figure 4.8 PP-attachment ambiguity

Coordination ambiguity occurs when it is not clear which phrases are being combined with a conjunction like *and*. For example, the phrase *beautiful hair and eyes* may have the structure [*beautiful hair*] and [*eyes*] or [*beautiful hair and*] and [*eyes*]. Identifying the correct parse from a number of possible parses is known as disambiguation. A parser may utilize statistical and semantic knowledge to disambiguate the parse tree, or it may return all possible parses and leave the disambiguation for subsequent processing. The basic top-down parser we discussed, returns the first successful parse without exploring other possibilities. It needs to be modified to return all possible parses.

A sentence may have *local ambiguity* resulting in inefficient parsing. Local ambiguity occurs when certain parts of a sentence are ambiguous. For example, the sentence *Paint the door* is unambiguous, but during parsing it is not known whether the first word *Paint* is a verb or a noun. Hence, the parser makes a few incorrect expansions before discovering that *Paint* is a verb. Thus, it must use backtracking, or parallelism, to consider both parses.

Yet another problem associated with our basic top-down strategy is that of repeated parsing. The parser often builds valid trees for portions of the input that it discards during backtracking. These have to be rebuilt during subsequent steps in the parse. If we could avoid this extra effort, we would have more efficient parsers.

Dynamic programming algorithms can solve these problems. These algorithms construct a table containing solutions to sub-problems, which, if solved, will solve the whole problem. In parsing, a dynamic programming

algorithm builds a table containing sub-trees for each and every constituent appearing in the input. The efficiency comes from the fact that once these sub-trees are discovered, they are stored and later consulted by all parses attempting to expand that constituent. This solves the re-parsing and ambiguity problem. There are three widely known dynamic parsers—the Cocke-YOUNGER-Kasami (CYK) algorithm, the Graham-Harrison-Ruzzo (GHR) algorithm, and the Earley algorithm. We now present the Earley and CYK parsing algorithms to illustrate top-down and a bottom-up dynamic programming algorithms respectively.

Probabilistic grammar can also be used to disambiguate parse trees. We discuss probabilistic parsing and how they can be used to identify a likely parse, in the next section.

4.4.4 Earley Parser

The Earley parser implements an efficient parallel top-down search using dynamic programming. It builds a table of sub-trees for each of the constituents in the input. This way, the algorithm eliminates the repetitive parse of a constituent which arises from backtracking, and successfully reduces the exponential-time problem to polynomial time. The Earley parser can handle recursive rules such as $A \rightarrow AC$ without getting into an infinite loop.

The most important component of this algorithm is the Earley chart that has $n+1$ entries, where n is the number of words in the input. The chart contains a set of states for each word position in the sentence. The algorithm makes a left to right scan of input to fill the elements in this chart. It builds a set of states, one for each position in the input string (starting from 0), that describe the condition of the recognition process at that point in the scan. The states in each entry provide the following information.

1. A sub-tree corresponding to a grammar rule.
2. Information about the progress made in completing the sub-tree.
3. Position of the subtree with respect to input.

A state is represented as a dotted rule and a pair of numbers representing starting position and the position of dot. This representation takes the form

$$A \rightarrow X_1 \dots \bullet C \dots X_m [i, j]$$

where the dot (\bullet) represents the position in the rule's right hand side, and the two numbers (i and j) represent where the state begins and where the dot lies. A dot at the right end of the rule represents a successful parse of the associated non-terminal.

Earley Parsing	Input: Sentence and the Grammar Output: Chart chart[0] $\leftarrow S' \rightarrow S, [0, 0]$ $n \leftarrow \text{length}(\text{sentence}) // \text{number of words in the sentence}$ for $i = 0$ to n do for each state in chart[i] do if [$\text{incomplete}(\text{state})$ and $\text{next category is a part of speech}$] then predictor [state] else if [$\text{incomplete}(\text{state})$ and $\text{next category is a part of speech}$] scanner [state] else completer [state] end-if end-for end-for Procedure predictor ($A \rightarrow X_1 \dots \bullet B \dots X_m [i, j]$) for each rule ($B \rightarrow \alpha$) in G do insert the state $B \rightarrow \bullet \alpha, [j, j]$ to chart[j] End Procedure scanner ($A \rightarrow X_1 \dots \bullet B \dots X_m [i, j]$) If B is one of the part of speech associated with word[j] then Insert the state $B \rightarrow \text{word}[j] \bullet, [i, j + 1]$ to chart[j + 1] End Procedure Completer ($A \rightarrow X_1 \dots \bullet, [i, k]$) for each $B \rightarrow X_1 \dots \bullet A \dots [i, j]$ in chart[j] do insert the state $B \rightarrow X_1 \dots A \bullet \dots [i, k]$ to chart[k] End
----------------	---

Figure 4.9 The Earley Parsing Algorithm

The algorithm uses three operations to process states in the chart. These are:

- Predictor
- Scanner
- Completer

The algorithm sequentially constructs the sets for each of the $n+1$ chart entries. Chart [0] is initialized with a dummy state $S' \rightarrow \bullet S, [0, 0]$. At each step one of the three operations are applicable depending on the state.

Completer

Application of these operators result in addition of new states to either the current or the next set of states. The presence of a state $S \rightarrow \alpha, [0,N]$ indicates a successful parse. The algorithm is shown in Figure 4.9.

We now explain the function of the three operators.

Predictor

As the name suggests, the predictor generates new states representing potential expansion of the non-terminal in the left-most derivation. A predictor is applied to every state that has a non-terminal to the right of the dot, when the category of that non-terminal is different from the part-of-speech. The application of this operator results in the creation of as many new states as there are grammar rules for the non-terminal. These new states are placed into the same chart entry as the generating state. Their start and end positions are at the point where the generating state ends. If

$$A \rightarrow X_1 \dots \bullet B \dots X_m [i, j]$$

Then for every rule of the form $B \rightarrow \alpha$, the operation adds to chart $[j]$, the state

$$B \rightarrow \bullet \alpha, [j, j]$$

For example, when the generating state is $S \rightarrow \bullet NP VP, [0,0]$, the predictor adds the following states to chart $[0]$:

$$\begin{aligned} NP &\rightarrow \bullet \text{ Det Nominal, } [0,0] \\ NP &\rightarrow \bullet \text{ Noun, } [0,0] \\ NP &\rightarrow \bullet \text{ Pronoun, } [0,0] \\ NP &\rightarrow \bullet \text{ Det Noun PP, } [0,0] \end{aligned}$$

Scanner

A scanner is used when a state has a part-of-speech category to the right of the dot. The scanner examines the input to see if the part-of-speech appearing to the right of the dot matches one of the part-of-speech associated with the current input. If yes, then it creates a new state using the rule that allows generation of the input word with this part-of-speech. It advances the pointer over the predicted input category and adds it to the next chart entry. If the state is $A \rightarrow \dots \bullet a, [i, j]$ and 'a' is one of the part-of-speech associated with w_j , then it adds $a \rightarrow \dots w_j \bullet, [i, j]$ to chart $[j+1]$.

Returning to our example, when the state $NP \rightarrow \bullet \text{ Det Nominal, } [0,0]$ is processed, the parser finds a part-of-speech category next to the dot. It checks if the category of the current word ($curr_word$) matches with the expectation in the current state. If yes, then it adds the new state $\text{Det} \rightarrow curr_word \bullet, [0,1]$ to the next chart entry.

Completer
The completer is used when the dot reaches the right end of the rule. The presence of such a state signifies successful completion of the parse of some grammatical category. The completer identifies all previously generated states that expect this grammatical category at this position in the input and creates new states by advancing the dots over the expected category. All these newly generated states are inserted in the current chart entry. More formally, if $A \rightarrow \dots \bullet, [i, k]$, then the completer adds $B \rightarrow \dots A \bullet \dots [i, k]$ to chart $[k]$ for all states $B \rightarrow \dots \bullet A \dots [i, j]$ in chart $[j]$. An item is added to a set only if it is not already in the set.

Example 4.1 Let us trace the algorithm using sentence (4.7). The sequence of states created by the parser is shown in Figure 4.10.

Chart [0]	S0	$S' \rightarrow \bullet S$	[0,0]
Dummy			
S1	$S \rightarrow \bullet NP VP$	start	
S2	$S \rightarrow \bullet VP$		
S3	$NP \rightarrow \bullet \text{ Det Nominal}$		
S4	$NP \rightarrow \bullet \text{ Pronoun}$		
S5	$NP \rightarrow \bullet \text{ Det Noun PP}$		
S6	$NP \rightarrow \bullet \text{ Verb NP}$		
S7	$VP \rightarrow \bullet \text{ Verb}$		
S8	$VP \rightarrow \bullet \text{ Verb NP}$		
S9	$\text{Noun} \rightarrow \text{paint} \bullet$	[0,1]	
S10	$\text{Verb} \rightarrow \text{paint} \bullet$	[0,1]	
S11	$NP \rightarrow \text{Noun} \bullet$	[0,1]	
S12	$VP \rightarrow \text{Verb} \bullet NP$	[0,1]	
S13	$VP \rightarrow \text{Verb} \bullet$	[0,1]	
S14	$S \rightarrow NP \bullet VP$	[0,1]	
S15	$NP \rightarrow \bullet \text{ Det Nominal}$	[1,1]	
S16	$NP \rightarrow \bullet \text{ Noun}$	[0,1]	
S17	$S \rightarrow VP \bullet$	[1,1]	
S18	$VP \rightarrow \bullet \text{ Verb NP}$	[1,1]	
S19	$VP \rightarrow \bullet \text{ Verb}$	[1,1]	
S20	$\text{Det} \rightarrow \text{the} \bullet$	[1,2]	
Chart [2]			
S21	$NP \rightarrow \text{Det} \bullet \text{ Nominal}$	[2,2]	
S22	$\text{Nominal} \rightarrow \bullet \text{ Noun}$	[2,2]	
S23	$\text{Nominal} \rightarrow \bullet \text{ Noun Nominal}$	[2,3]	
S24	$\text{Noun} \rightarrow \text{door} \bullet$	[2,3]	
Chart [3]			
S25	$Nominal \rightarrow \bullet \text{ Noun} \bullet$	[1,3]	
S26	$NP \rightarrow \text{Det Nominal} \bullet$	[0,3]	
S27	$S \rightarrow NP \bullet VP$	[0,3]	
S28	$VP \rightarrow \text{Verb NP} \bullet$	[3,3]	
S29	$VP \rightarrow \bullet \text{ Verb NP}$	[3,3]	
S30	$VP \rightarrow \bullet \text{ Verb}$	[0,3]	
S31	$S \rightarrow VP \bullet$		

Figure 4.10 Sequence of states created in parsing sentence (4.7) using Earley algorithm

The presence of the state $S \rightarrow VP \bullet, [0,3]$ in the chart indicates successful completion of the parse of the sentence, but it does not return the exact parse. Thus, the algorithm in this form can be used only to recognize a sentence. However, it is possible to utilize the entries appearing in the table to construct the parse tree of a valid sentence. In order to do so the algorithm needs to be modified. This modification is carried out by the completer which creates new states by advancing earlier incomplete states whenever it finds a state having a dot at the end of the rule. We add a pointer to the previous states of the new state. A list of previous states is thus maintained. Extracting a parse tree from the chart involves following these pointers, beginning with the state marking the successful completion of the parse. Though the Earley algorithm fills the chart entry in polynomial time, extracting all parse trees still requires an exponential amount of time.

The Earley parser can be augmented with unification structures to eliminate ill-formed structures as they are introduced. This requires certain modification in the algorithm. The first change involves the addition of a feature structure derived from their unification constraints. The second change is the addition of a new field to the chart. This new field is a directed acyclic graph representing the feature structure corresponding to the state. The details of the modified Earley algorithm can be found in Jurafsky and Martin (2000).

4.4.5 The CYK Parser

Like the Earley algorithm, the CYK (Cocke-YOUNGER-Kasami) is a dynamic programming parsing algorithm. However, it follows a bottom-up approach in parsing. It builds a parse tree incrementally. Each entry in the table is based on previous entries. The process is iterated until the entire sentence has been parsed. The CYK parsing algorithm assumes the grammar to be in Chomsky normal form (CNF). A CFG is in CNF if all the rules are of only two forms:

$$A \rightarrow BC$$

$A \rightarrow w$, where w is a word.

The algorithm first builds parse trees of length one by considering all rules which could produce words in the sentence being parsed. Then, it constructs the most probable parse for all the constituents of length two. The parse of shorter constituents constructed in earlier iterations can be used in constructing the parse of longer constituents.

Like the Earley algorithm, the basic CYK algorithm is also a chart-based algorithm. A non-terminal is stored in the $[i, j]$ th entry of the chart if, and only if, $A \Rightarrow w_i \cdot w_{i+1} \dots w_{i+j-1}$. The chart is triangular. A sentence is recognized if the start symbol S is in the entry $[1, n]$ of the chart. Beginning with the start symbol of the grammar, we are able to derive the entire sequence of words appearing in the sentence. The algorithm builds smaller constituents before attempting to construct larger ones. First, the terminal derivation rules of the grammar are used to generate the $[i, 1]$ th entries. These entries represent non-terminals that derive the individual words appearing in the sentence, $w_{i1} = w_p$, for $1 \leq i \leq n$, where n is the length (number of words) of the sentence. $A \Rightarrow w_{i1}$ if $A \rightarrow w_i$ is a rule in the grammar. It then continues with sub-string of length two, three, and so on. For every non-terminal A in the grammar, the algorithm determines if $A^* \Rightarrow w_{ij}^*$. As the grammar is in CNF, A could derive if there existed a rule of the form $A \rightarrow BC$ such that B derives the first k words of the w_{ij} (i.e. $B \rightarrow w_{ik}$) and C derives the remaining $j-k$ words ($C \rightarrow w_{ij}^*$) as shown in Figure 4.11. More formally,

$$A^* \Rightarrow w_{ij} \text{ if}$$

1. $A \rightarrow BC$ is a rule in grammar
2. $B^* \Rightarrow w_{ik}$ and
3. $C^* \Rightarrow w_{ij}^*$

For a sub-string w_{ij} of length j starting at i , the algorithm considers all possible ways of breaking it into two parts w_{ik} and w_{ij}^* . Finally, since $s = w_{1n}$ we have to verify that $S^* \Rightarrow w_{1n}$ i.e., the start symbol of the grammar derives w_{1n} .

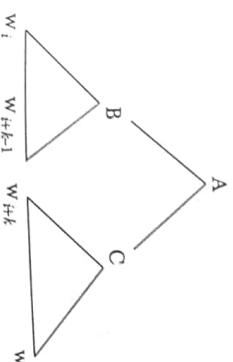


Figure 4.11 Breaking a string

The steps involved in the algorithm are shown in Figure 4.12. In order to use the information contained in the table to construct parse trees, we need to maintain back-pointers to the table entries we combine. This allows us to construct all possible parse trees by following back pointers.

4.5 PROBABILISTIC PARSING

```

Let  $w = w_1 w_2 w_3 \dots w_j \dots w_n$ 
and  $w_{ij} = w_i \dots w_{i+j-1}$ 
// Initialization step
for  $i := 1$  to  $n$  do
  for all rules  $A \rightarrow w_i$  do
    chart [ $i, i$ ] = { $A$ }
  // Recursive step
  for  $j = 2$  to  $n$  do
    for  $i = 1$  to  $n-j+1$  do
      begin
        chart [ $i, j$ ] =  $\emptyset$ 
        for  $k = 1$  to  $j-1$  do
          chart [ $i, j$ ] := chart [ $i, j$ ]  $\cup$  { $A \mid A \rightarrow BC$  is a production and
           $B \in \text{chart}[i, k]$  and  $C \in \text{chart}[i+k, j-k]$ }
        end
        if  $S \in \text{chart}[1, n]$  then accept else reject
      end
    end
  end
}

```

Figure 4.12 The CYK algorithm

To give a better understanding of the whole idea, we work out an example. Consider the following simplified grammar in CNF:

$$\begin{array}{ll}
 S \rightarrow NP\ VP & \\
 VP \rightarrow Verb\ NP & \quad \text{Verb} \rightarrow \text{wrote} \\
 NP \rightarrow Det\ Noun & \quad \text{Noun} \rightarrow \text{girl} \\
 \text{Def} \rightarrow \text{an} \mid \text{the} & \quad \text{Noun} \rightarrow \text{essay}
 \end{array}$$

The sentence to be parsed is: *The girl wrote an essay*.

Table 4.4 contains entries after a complete scan of the algorithm. The entry in the $[1, n]$ th cell contains a start symbol which indicates that $S^* \Rightarrow w_1 \dots$, i.e., the parse is successful. It is possible for a cell to have multiple entries.

Table 4.4 Sequence of states created in the chart by the CYK algorithm while parsing the sentence, *Sara wrote an essay*

1	2	3	4	5
Det \rightarrow The	NP \rightarrow Det Noun			
Noun \rightarrow Girl				
Verb \rightarrow wrote		VP \rightarrow Verb NP		
Det \rightarrow an	NP \rightarrow Det Noun			
Noun \rightarrow essay				

In order to construct a statistical parser, we have to first find all possible parses of a sentence, then assign probabilities to them, and finally return the most probable parse. We discuss an implementation of statistical parsing based upon probabilistic context-free grammars (PCFGs). But before proceeding ahead to this discussion, let us have a look at why we need probabilistic parsing at all. What advantages do these parsers offer? The first benefit that a probabilistic parser offers is removal of ambiguity from parsing. We have seen earlier, that sentences can have multiple parse trees. The parsing algorithm discussed so far in this chapter has no means to decide which parse is the correct or most appropriate one. Probabilistic parsers assign probabilities to parses. These probabilities are then used to decide the most likely parse tree structure of an input sentence. Another benefit this parser offers is related to efficiency. The search space of possible tree structures is usually very large. With no information on which sub-trees are more likely to be a part of the final parse tree, the search can be quite time consuming. Using probabilities to guide the process, the search becomes more efficient.

A probabilistic context-free grammar (PCFG) is a CFG in which every rule is assigned a probability (Charniak 1993). It extends the CFG by augmenting each rule $A \rightarrow \alpha$ in set of productions P , with a conditional probability p :

$$A \rightarrow \alpha [p]$$

where p gives the probability of expanding a constituent using the rule $A \rightarrow \alpha$.

Let us now define PCFG. A PCFG is defined by the pair (G, f) , where G is a CFG and f is a positive function defined over the set of rules such that, the sum of the probabilities associated with the rules expanding a particular non-terminal is 1 (Infante-Lopez and Maarten de Rijke 2006).

$$\sum_{\alpha} f(A \rightarrow \alpha) = 1$$

An example of PCFG is shown in Figure 4.13. We can verify that for each non-terminal, the sum of probabilities is 1.

$$\begin{aligned} f(S \rightarrow NP VP) + f(S \rightarrow VP) &= 1 \\ f(NP \rightarrow Det Noun) + f(NP \rightarrow Noun) + f(NP \rightarrow Pronoun) + f(NP \\ \rightarrow Det Noun PP) &= 1 \\ f(VP \rightarrow Verb NP) + f(NP \rightarrow Verb) + f(VP \rightarrow VP PP) &= 1.0 \\ f(Det \rightarrow this) + f(Det \rightarrow that) + f(Det \rightarrow a) + f(Det \rightarrow the) &= 1.0 \\ f(Noun \rightarrow paint) + f(Noun \rightarrow door) + f(Noun \rightarrow bird) + f(Noun \\ \rightarrow hole) &= 1.0 \end{aligned}$$

Rule	Count $(A \rightarrow \alpha)$	Count A	MLE estimates
$S \rightarrow VP$	2	2	1
$NP \rightarrow Det Noun PP$	1	4	0.25
$NP \rightarrow Det Noun$	3	4	0.75
$VP \rightarrow Verb NP$	2	3	0.66
$VP \rightarrow VP PP$	1	3	0.33
$Det \rightarrow the$	2	2	1
$Noun \rightarrow hole$	2	4	0.5
$Noun \rightarrow door$	2	4	0.5
$Prep \rightarrow with$	1	1	1
$Verb \rightarrow Paint$	1	1	1

Figure 4.13 A probabilistic context-free grammar (PCFG)

Similarly, the condition is satisfied for the other non-terminals.

4.5.1 Estimating Rule Probabilities

The next question is how are probabilities assigned to rules? One way to estimate probabilities for a PCFG is to manually construct a corpus of a parse tree for a set of sentences, and then estimate the probabilities of each rule being used by counting them over the corpus. The MLE estimate for a rule $A \rightarrow \alpha$ is given by the expression

$$P_{MLE}(A \rightarrow \alpha) = \frac{\text{Count}(A \rightarrow \alpha)}{\sum_{\alpha} \text{Count}(A \rightarrow \alpha)}$$

If our training corpus consists of two parse trees (as shown in Figure 4.14), we will get the estimates as shown in Table 4.5 for the rules:

Table 4.5 The MLE for the grammar rules used in trees of Figure 4.14

Rule	Count $(A \rightarrow \alpha)$	Count A	MLE estimates
$S \rightarrow VP$	2	2	1
$NP \rightarrow Det Noun PP$	1	4	0.25
$NP \rightarrow Det Noun$	3	4	0.75
$VP \rightarrow Verb NP$	2	3	0.66
$VP \rightarrow VP PP$	1	3	0.33
$Det \rightarrow the$	2	2	1
$Noun \rightarrow hole$	2	4	0.5
$Noun \rightarrow door$	2	4	0.5
$Prep \rightarrow with$	1	1	1
$Verb \rightarrow Paint$	1	1	1

We now turn to another important question—what do we do with these probabilities? We assign a probability to each parse tree φ of a sentence s . The probability of a complete parse is calculated by multiplying the probabilities for each of the rules used in generating the parse tree:

$$P(\varphi, s) = \prod_{n \in \varphi} p\{r(n)\}$$

where n is a node in the parse tree φ and r is the rule used to expand n .

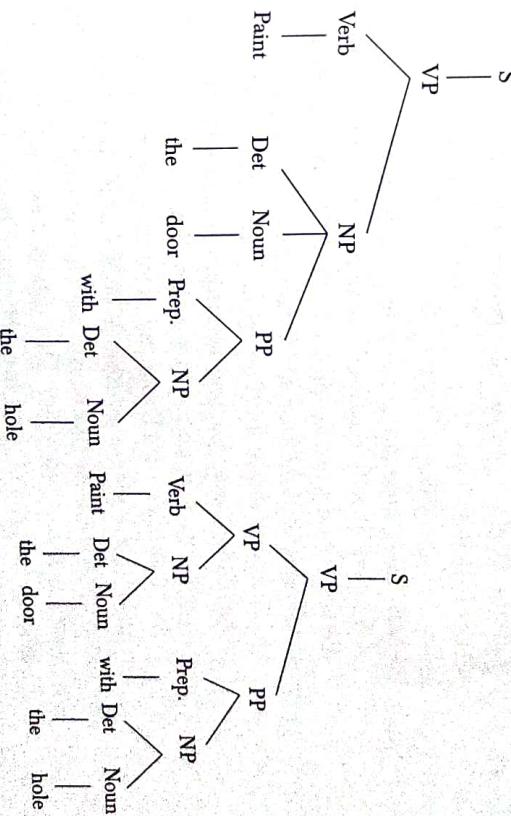


Figure 4.14 Two possible parse trees

For a sentence, the parse trees generated by a PCFG are the same as those generated by a corresponding CFG. However, the PCFG assigns a probability to each parse. The probability of the two parse trees of the sentence *Paint the door with the hole* (shown in Figure 4.14) can be computed as follows:

$$\begin{aligned} P(t1) &= 0.2 * 0.5 * 0.2 * 0.2 * 0.35 * 0.25 * 1.0 * 0.25 * 0.4 * 0.35 * 0.25 \\ &= 0.0000030625 \\ P(t2) &= 0.2 * 0.2 * 0.5 * 0.2 * 0.4 * 0.35 * 0.25 * 1 * 0.25 * 0.4 * 0.35 * 0.25 \\ &= 0.000001225 \end{aligned}$$

In this example, the first tree has a higher probability leading to correct interpretation.

We can calculate probability to a sentence s by summing up probabilities of all possible parses associated with it.

$$P(s) = \sum_{t \in \tau(s)} P(t, s)$$

Thus, the sentence will have the probability

$$\begin{aligned} P(t1) + P(t2) &= 0.0000030625 + 0.000001225 \\ &= 0.0000042875 \end{aligned}$$

4.5.2 Parsing PCFGs

Given a PCFG, a probabilistic parsing algorithm assigns the most likely parse ϕ' to a sentence s .

$$\phi' = \operatorname{argmax}_{T \in \tau(S)} P(T | S)$$

where $\tau(S)$ is the set of all possible parse trees of S .

We have already discussed a number of parsing algorithms for CFG. We will now discuss the probabilistic CYK parsing algorithm. We have already discussed the basic CYK parsing algorithm. Its probabilistic version was first introduced by Ney (1991) and is generally preferred over probabilistic Earley Parsing algorithm due to its simplicity. We begin with notations that we use in the algorithm.

As in the basic CYK Parsing algorithm, $w = w_1 w_2 w_3 w_4 \dots w_n$ represents a sentence consisting of n words. Let $\phi[i, j, A]$ represent the maximum probability parse for a constituent with non-terminal A spanning words $i, i+1, \dots, j-1$. This means it is a sub-tree rooted at A that derives sequence of j words beginning at position i and has a probability greater than all other possible sub-trees deriving the same word sequence.

```

Initialization:
  for  $i := 1$  to  $n$  do
    for all rules  $A \rightarrow w_i$  do
       $\phi[i, 1, A] = R(A \rightarrow w_i)$ 

Recursive Step:
  for  $j = 2$  to  $n$  do
    for  $i = 1$  to  $n-j+1$  do
      begin
         $\phi[i, 1, A] = \phi$ 
        for  $k = 1$  to  $j-1$  do
           $\phi'[i, j, A] = \max_k \phi'[i, k, B] \times \phi'[k, j, C] \times R(A \rightarrow BC),$ 
          such that  $A \rightarrow BC$  is a production rule in grammar
           $BP[i, j, A] = \{k, A, B\}$ 
      end

```

Figure 4.15 Probabilistic CYK algorithm

The algorithm is given in Figure 4.15. Like the basic CYK algorithm, the first step is to generate items of length 1. However, in this case, we have to initialize the maximum probable parse trees deriving a string of length 1, with the probabilities of the terminal derivation rules used to derive them.

$\phi'[i, 1, A] = P(A \rightarrow w_i)$

Similarly, recursive step involves breaking a string into all possible ways and identifying the maximum probable parse.

$$\phi'[i, j, A] = \max_k \phi'[i, k, B] \times \phi'[k, j, C] \times R(A \rightarrow BC)$$

The rest of the steps follow those of basic CYK parsing algorithm.

4.5.3 Problems with PCFG

The PCFG is not without disadvantages. Its first problem lies in the independence assumption. We calculate the probability of a parse tree assuming that the rules are independent of each other. However, this is not true. How a node expands depends on its location in the parse tree. For example, Francis et al. (1999) showed that pronouns occur more

frequently as subjects rather than objects. These dependencies are not captured by a PCFG, as the probability of, say, expanding an NP as a pronoun versus a lexical NP, is independent of whether the NP appears as a subject or an object.

Another problem associated with a PCFG is its lack of sensitivity to lexical information. Lexical information plays a major role in determining correct parse in case of PP attachment ambiguities and coordination ambiguities (Collins 1999). Two structurally different parses that use the same rules will have the same probability under a PCFG, making it difficult to identify the correct or most probable parse. The words appearing in a parse may make certain parses unnatural. This however, requires a model which captures lexical dependency statistics for different words. Such a model is presented next.

Lexicalization

In PCFG, the chance of a non-terminal expanding using a particular rule is independent of the actual words involved. However, this independence assumption does not seem reasonable. Words do affect the choice of the rule. Investigations made on tree bank data suggest that the probabilities of various common sub-categorization frames differ depending on the verb that heads the verb phrase (Manning and Schütze 1999). This suggests the need for lexicalization, i.e., involvement of actual words in the sentences, to decide the structure of the parse tree. Lexicalization is also helpful in choosing phrasal attachment positions. This model of lexicalization is based on the idea that these are strong lexical dependencies between heads and their dependents, for example, between a head noun and its modifiers, or between a verb and a noun phrase object, where the noun phrase object in turn can be approximated by its head noun.

One way to achieve lexicalization is to mark each phrasal node in a parse tree by its head word. Figure 4.16 is an example of such a tree. One way to implement this model is to use a lexicalized grammar. A lexicalized grammar is a grammar in which every finite structure is associated with one or more lexical heads. A context free grammar is not lexicalized as no lexical item is associated with its rule, such as, $S \rightarrow NP VP$. Nor is the PCFG lexicalized. In order to convert a PCFG into lexicalized grammar, each of its rules must identify a head daughter, which is one of the constituents appearing on its right hand side, for example head daughter of $S \rightarrow NP VP$ rule is VP . The head word of a node in the parse tree is set to the head word of its head daughter. For example, the head of a verb phrase is the main verb. Hence, the head of the node VP in the Figure 4.16 is jumped. Similarly, the head of a constituent expanded using the rule $S \rightarrow NP VP$ is the head of VP .

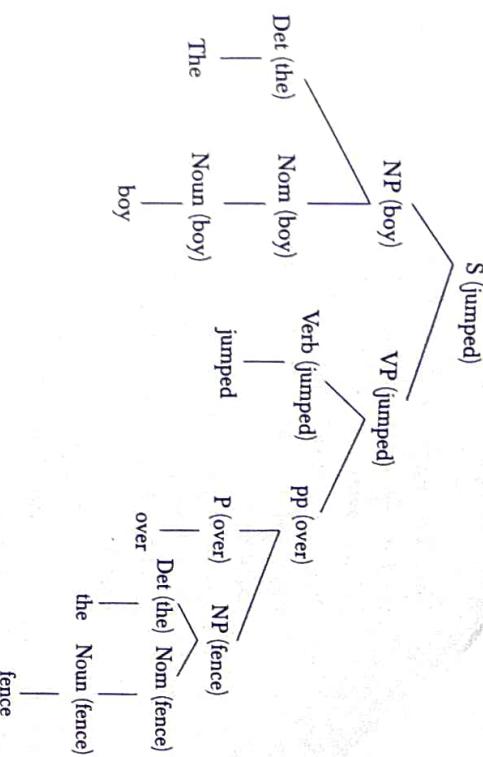


Figure 4.16 A lexicalized tree

The probability in a lexicalized PCFG is calculated for each rule or head of a phrase and head of the sub-phrase or head of the phrase (i.e., head/mother head) combination. For example, we need to collect the following rule/head probability:

$$\begin{aligned} VP(\text{jumped}) &\rightarrow \text{verb } (\text{jumped}) \text{ PP } (\text{over}) \\ VP(\text{jumped}) &\rightarrow \text{verb } (\text{jumped}) \text{ PP } (\text{into}) \\ VP(\text{jumped}) &\rightarrow \text{verb } (\text{jumped}) \text{ PP } (\text{to}) \end{aligned}$$

An example of the head or mother head combination is the following type of probability:

"What is the probability that an NP whose mother's head is *over* has the head *fence*?"

With a lexicalized PCFG, the probability of the parse tree is computed as the product of the probability of the head of the constituent c given the probability of the mother $m(c)$ and the probability of the rule used to expand constituent c given the head of constituent (Charniak 1997):

$$P(\varphi, s) = \prod_{c \in \varphi} p[(h(c)|m(c)) \cdot p[(r(c)|h(c))]$$

where $h(c)$ = head of the constituent c

$r(c)$ = rule used to expand constituent c

and $m(c)$ = mother of the constituent c

Not all natural languages have the same characteristics. So far we have talked only of the English language, using CFG as the grammar formalism.

But CFG may not be a viable choice for other languages. This is particularly true for Indian languages. In this section, we discuss some of the characteristics of Indian languages that make CFG unsuitable. We then give a brief outline of a parser-based on Paninian grammar. As discussed in Chapter 2, Paninian grammar can be used to model Indian languages.

The majority of the Indian languages are free word order. By free word order language we mean that the order of words in a sentence can be changed without leading to a grammatically incorrect sentence. For example:

सबा खाना खाती है । Saba khana khati hai.

खाना सबा खाती है । Khana Saba khati hai.

Both are valid Hindi sentences meaning *Saba eats food*.

The same is true for Urdu and most other Indian languages. The CFG we used for parsing English is basically positional. It can be used to model language in which a position of the constituents carries useful information, but it fails to model free word order languages.

Extensive and productive use of complex predicates (CPs) is another property that most Indian languages have in common. A complex predicate combines a light verb with a verb, noun, or adjective, to produce a new verb. For example,

(a) सबा आयी ।

(*Saba Ayi*)

Saba came.

(b) सबा आ गयी ।

(*Saba a gayi*)

Saba came went.

(c) सबा आ पड़ी ।

Saba a pari.

Saba come fell.

Saba came (suddenly).

The complex predicates change the functional structure of the sentence; however they do not change the sub-categorization frame of the verb.

The use of post-position case markers (Vibhakti) and verb complexes consisting of sequences of verbs, e.g., खा रही है (*kha rahi hai*) are other properties common to Indian languages. The auxiliary verbs in this grammar provides a framework to model Indian languages. It focuses on the extraction of Karak relations from a sentence.

Bharti and Sangal (1990) described an approach for parsing of Indian languages based on Paninian grammar formalism. Their parser works in two stages. The first stage is responsible for identifying word groups in the second for assigning a parse structure to the input sentence. The input to the first stage comes from the morphological analysis phases (as discussed in Chapter 3). The output of the first stage is word groups, which are sequences of words that act as a unit. For example, a verb group consists of a main verb and a sequence of auxiliaries. For the sentence

लड़कियाँ मैदान में हॉकी क्रील रही हैं ।

Ladkiyan maidaan mein hockey khel rahi hein.

The word *ladkiyan* forms one unit, the words *maidaan* and *mein* are grouped together to form a noun group, and the word sequence *khel rahi* *hein* forms a verb group.

The choice of noun and verb groups over noun and verb phrases adds computational simplicity to the approach. The concept of verb phrase is not natural to Indian languages and computing a noun phrase is difficult.

In the second stage, the parser takes the word groups formed during first stage and identifies (i) Karaka relations among them, and (ii) senses of words. A data structure, called Karaka chart, stores additional information like Karaka-Vibhakti mapping, Karaka necessity (mandatory or optional), and transformation rules for Karaka relations, needed in this step. Transformation rules tell us how to create a Karaka chart for a verb group using the default Karaka chart. The form of the default Karaka chart is shown in Table 4.6.

Table 4.6 Default Karaka chart

Karaka (case relations)	Vibhakti (case markers or post-positions)	Necessity
Karta	φ	Mandatory
Karma	Ko or φ	Mandatory
Adhikaran	Mein or par	Optional
Sampradan	Ko or ke liye	Optional

Once the Karaka chart for the verb groups are ready, noun groups are tested against them. A noun group satisfying the Vibhakti restriction for a verb group becomes a candidate for its Karaka. The Karaka relation between a verb group and a noun group can be depicted using a constraint graph. Nodes in the graph represent word groups and an arc represents a Karaka restriction between word groups. We have earlier identified the word groups in sentence (4.8). Its constraint graph is given in Figure 4.17.

SUMMARY

- Syntactic parsing deals with the syntactic structure of a sentence.
- In many languages, words are grouped to form larger groups termed constituent or phrases, which can be modelled by context-free grammar (CFG).
- A CFG consists of a set of rules or productions stating which elements can occur in a phrase and in which order.

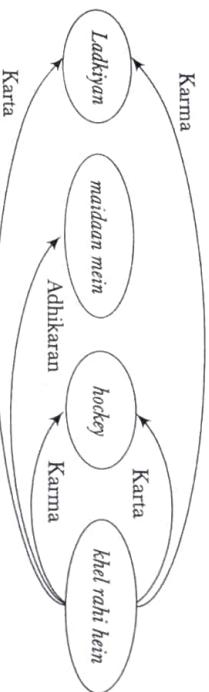


Figure 4.17 Constraint graph for sentence (4.8)

Each sub-graph of the constraint graph that satisfies the following constraints yields a parse of the sentence.

1. It contains all the nodes of the graph.
2. It contains exactly one outgoing edge from a verb group for each of its mandatory Karakas. These edges are labelled by the corresponding Karaka.
3. For each of the optional Karaka in Karaka chart, the sub-graph can have at most one outgoing edge labelled by the Karaka from the verb group.
4. For each noun group, the sub-graph should have exactly one incoming edge.

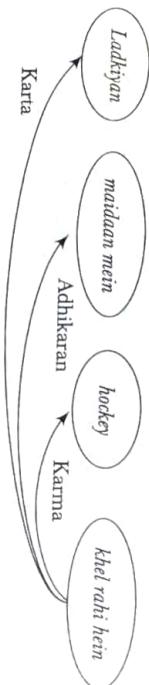


Figure 4.18 A parse of the sentence (4.8)

A sub-graph of the constraint graph (Figure 4.17) satisfying these constraints is shown in Figure 4.18. This represents one of the possible parse of sentence (4.8).

More than one sub-graph may satisfy these constraints in which case the sentence is ambiguous and has multiple parse associated with it. If no sub-graph satisfies these constraints, then the grammar fails to assign any parse structure to the sentence. Disambiguation among various senses of verbs and nouns is carried out with the help of the Lakshman chart, which contains features that are useful for discriminating among various senses. Readers are referred to Bharti and Sangal (1990) and Bharti et al. (1995) for detailed treatment of the algorithm.

REFERENCES

- Backus, J.W., 1959, 'The syntax and semantics of the proposed international algebraic language of the Zurich ACM-GAMM conference,' *Proceedings of the International Conference on Information Processing*, UNESCO, pp. 125-32.

- Bharti, Akshar and Rajeev Sangal, 1990, 'A Karaka-based approach to parsing of Indian languages,' *Proceedings of the 13th Conference on Computational Linguistics*, Association for Computational Linguistics, 3.
- Bharti, Akshar, Vineet Chaitanya, and Rajeev Sangal, 1995, *Natural Language Processing: A Paninian Perspective*, Prentice-Hall of India.
- Charniak, Eugene, 1993, *Statistical Language Learning*, MIT Press, Cambridge.
- , 1997, 'Statistical techniques for natural language parsing,' *AI Magazine*.
- Chomsky, N., 1957, *Syntactic Structures*, Mouton, The Hague.
- Collins, M.J., 'Head-driven statistical parsing for natural language processing,' *Ph.D. Thesis*, University of Pennsylvania, Philadelphia.
- Infante-Lopez, Gabriel and Maarten de Rijke, 2006, 'A note on the expressive power of probabilistic context free grammars,' *Journal of Logic, Language and Information*, Kluwer Academic Publisher, 15(3).
- Jurafsky, Daniel and James H. Martin, 2000, *Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics and Speech Recognition*, Prentice Hall, NJ.
- Manning, C. and H. Schütze, 1999, *Foundations of Statistical Natural Language Processing*, MIT Press, Cambridge.
- Marcus, Mitchell P., Beatrice Santorini, and Mary Ann Marcinkiewicz, 1993, 'Building a large annotated corpus of English: the Penn treebank,' *Computational Linguistics*, 19, pp. 313–30.
- Naur, Peter, J.W. Backus, F.L. Bauer, J. Green, C. Katz, J. McCarthy, A. J. Perlis, H. Rutishauser, K. Samuelson, B. Vauquois, J. H. Wegstein, A. van Wijngaarden, and M. Woodger, 1960, 'Report on the algorithmic language ALGOL 60,' *Communications of the ACM*, 3(5), pp. 299–314.
- Ney, H., 1991, 'Dynamic programming parsing for context-free grammars in continuous speech recognition,' *IEEE Transactions on Signal Processing*, 39(2), pp. 336–40.

EXERCISES

- Give two possible parse trees for the sentence, *Stolen painting found by tree*.
- Identify the noun and verb phrases in the sentence, *My soul answers in music*.
- Give the correct parse of sentence (4.6).

LAB EXERCISES

- Discuss the disadvantages of the basic top-down parser with the help of an appropriate example.
- Tabulate the sequence of states created by CYK algorithm while parsing, *The sun rises in the east*. Augment the grammar in section 4.4.5 with appropriate rules of lexicon.
- Discuss the disadvantages of probabilistic context free grammar.
- What does lexicalized grammar mean? How can lexicalization be achieved? Explain with the help of suitable examples.
- List the characteristics of a garden path sentence. Give an example of a garden path sentence and show its correct parse.
- What is the need of lexicalization?
- Use the following grammar:

$S \rightarrow NP VP$	$S \rightarrow VP$	$NP \rightarrow Det Noun$
$NP \rightarrow Noun$	$NP \rightarrow NP PP$	$VP \rightarrow VP NP$
$VP \rightarrow Verb$	$VP \rightarrow VP PP$	$PP \rightarrow Preposition NP$

Give two possible parse of the sentence: '*Pluck the flower with the stick*'. Introduce lexicon rules for words appearing in the sentence. Using these parse trees obtain maximum likelihood estimates for the grammar rules used in the tree. Calculate probability of any one parse tree using these estimates.