

Last Class: Smoothing

Today: Parsing Intro

1. Grammars and parsing
2. Top-down and bottom-up parsing

Slide CS474-1

Syntax

syntax: from the Greek *syntaxis*, meaning “setting out together or arrangement.”

Refers to the way words are arranged together.

Why worry about syntax?

- The boy ate the frog.
- The frog was eaten by the boy.
- The frog that the boy ate died.
- The boy whom the frog was eaten by died.

Slide CS474-2

Syntactic Analysis

Key ideas:

- **constituency:** groups of words may behave as a single unit or phrase
- **grammatical relations:** refer to the SUBJECT, OBJECT, INDIRECT OBJECT, etc.
- **subcategorization and dependencies:** refer to certain kinds of relations between words and phrases, e.g. *want* can be followed by an infinitive, but *find* and *work* cannot.

All can be modeled by various kinds of grammars that are based on context-free grammars.

Slide CS474-3

Grammars and Parsing

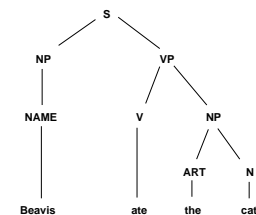
Need a **grammar:** a formal specification of the structures allowable in the language.

Need a **parser:** algorithm for assigning syntactic structure to an input sentence.

Sentence

Beavis ate the cat.

Parse Tree



Slide CS474-4

CFG example

CFG's are also called phrase-structure grammars.

Equivalent to Backus-Naur Form (BNF).

- | | |
|---------------------------|-------------------------------------|
| 1. $S \rightarrow NP VP$ | 5. $NAME \rightarrow \text{Beavis}$ |
| 2. $VP \rightarrow V NP$ | 6. $V \rightarrow \text{ate}$ |
| 3. $NP \rightarrow NAME$ | 7. $ART \rightarrow \text{the}$ |
| 4. $NP \rightarrow ART N$ | 8. $N \rightarrow \text{cat}$ |
- CFG's are *powerful* enough to describe most of the structure in natural languages.
 - CFG's are *restricted* enough so that efficient parsers can be built.

Slide CS474-5

CFG's

A context free grammar consists of:

1. a set of non-terminal symbols N
2. a set of terminal symbols Σ (disjoint from N)
3. a set of productions, P , each of the form $A \rightarrow \alpha$, where A is a non-terminal and α is a string of symbols from the infinite set of strings $(\Sigma \cup N)^*$
4. a designated start symbol S

Slide CS474-6

Derivations

- If the rule $A \rightarrow \beta \in P$, and α and γ are strings in the set $(\Sigma \cup N)^*$, then we say that $\alpha A \gamma$ **directly derives** $\alpha \beta \gamma$, or $\alpha A \gamma \Rightarrow \alpha \beta \gamma$
- Let $\alpha_1, \alpha_2, \dots, \alpha_m$ be strings in $(\Sigma \cup N)^*$, $m > 1$, such that

$$\alpha_1 \Rightarrow \alpha_2, \alpha_2 \Rightarrow \alpha_3, \dots, \alpha_{m-1} \Rightarrow \alpha_m,$$

then we say that α_1 **derives** α_m or $\alpha_1 \xRightarrow{*} \alpha_m$

Slide CS474-7

L_G

The language L_G generated by a grammar G is the set of strings composed of terminal symbols that can be derived from the designated start symbol S .

$$L_G = \{w | w \in \Sigma^*, S \xRightarrow{*} w\}$$

Parsing: the problem of mapping from a string of words to its parse tree according to a grammar G .

Slide CS474-8

General Parsing Strategies

Grammar	Top-Down	Bottom-Up
1. $S \rightarrow NP VP$	$S \rightarrow NP VP$	$\rightarrow \text{NAME ate the cat}$
2. $VP \rightarrow V NP$	$\rightarrow \text{NAME VP}$	$\rightarrow \text{NAME V the cat}$
3. $NP \rightarrow \text{NAME}$	$\rightarrow \text{Beav VP}$	$\rightarrow \text{NAME V ART cat}$
4. $NP \rightarrow \text{ART N}$	$\rightarrow \text{Beav V NP}$	$\rightarrow \text{NAME V ART N}$
5. $\text{NAME} \rightarrow \text{Beavis}$	$\rightarrow \text{Beav ate NP}$	$\rightarrow \text{NP V ART N}$
6. $V \rightarrow \text{ate}$	$\rightarrow \text{Beav ate ART N}$	$\rightarrow \text{NP V NP}$
7. $\text{ART} \rightarrow \text{the}$	$\rightarrow \text{Beav ate the N}$	$\rightarrow \text{NP VP}$
8. $N \rightarrow \text{cat}$	$\rightarrow \text{Beav ate the cat}$	$\rightarrow S$

Slide CS474–9

A Top-Down Parser

Input: CFG grammar, lexicon, sentence to parse

Output: yes/no

State of the parse: (*symbol list, position*)

₁ The ₂ old ₃ man ₄ cried ₅

start state: ((S) 1)

Slide CS474–10

Grammar and Lexicon

Grammar:

- $S \rightarrow NP VP$
- $NP \rightarrow \text{art } n$
- $NP \rightarrow \text{art adj } n$

- $VP \rightarrow v$
- $VP \rightarrow v NP$

Lexicon:

the: art
old: adj, n
man: n, v
cried: v

₁ The ₂ old ₃ man ₄ cried ₅

Slide CS474–11

Algorithm for a Top-Down Parser

$PSL \leftarrow (((S) 1))$

- Check for failure.* If PSL is empty, return NO.
- Select the current state, C.* $C \leftarrow \text{pop}(PSL)$.
- Check for success.* If $C = (()) <\text{final-position}>$, YES.
- Otherwise, generate the next possible states.*
 - $s_1 \leftarrow \text{first-symbol}(C)$
 - If s_1 is a *lexical symbol* and next word can be in that class, create new state by removing s_1 , updating the word position, and adding it to PSL. (I'll add to front.)
 - If s_1 is a *non-terminal*, generate a new state for each rule in the grammar that can rewrite s_1 . Add all to PSL. (Add to front.)

Slide CS474–12

Example

Current state

1. ((S) 1)
2. ((NP VP) 1)
3. ((art n VP) 1)
4. ((n VP) 2)
5. ((VP) 3)
6. ((v) 3)
7. (() 4)

Backup states

- ((art adj n VP) 1)
- ((art adj n VP) 1)
- ((art adj n VP) 1)
- ((v NP) 3) ((art adj n VP) 1)
- ((v NP) 3) ((art adj n VP) 1) Backtrack

Slide CS474–13

- | | | |
|-----------------------|--|--|
| 8. ((v NP) 3) | | ((art adj n VP) 1) leads to backtracking |
| ... | | |
| 9. ((art adj n VP) 1) | | |
| 10. ((adj n VP) 2) | | |
| 11. ((n VP) 3) | | |
| 12. ((VP) 4) | | |
| 13. ((v) 4) | | ((v NP) 4) |
| 14. (() 5) | | ((v NP) 4) |
| YES | | |

DONE!

Slide CS474–14

Problems with the Top-Down Parser

1. Only judges grammaticality.
2. Stops when it finds a single derivation.
3. No semantic knowledge employed.
4. No way to rank the derivations.
5. Problems with left-recursive rules.
6. Problems with ungrammatical sentences.

Slide CS474–15

Efficient Parsing

The top-down parser is terribly inefficient.

Have the first year Phd students in the computer science department take the Q-exam.

Have the first year Phd students in the computer science department taken the Q-exam?

Slide CS474–16