

Automata and Formal Languages (7)

Email: christian.urban at kcl.ac.uk
Office: SI.27 (1st floor Strand Building)
Slides: KEATS (also home work is there)

Two Weeks Ago: CFGs

A **context-free** grammar (CFG) G consists of:

- a finite set of nonterminal symbols (upper case)
- a finite terminal symbols or tokens (lower case)
- a start symbol (which must be a nonterminal)
- a set of rules

$$A \rightarrow \text{rhs}_1 | \text{rhs}_2 | \dots$$

where **rhs** are sequences involving terminals and nonterminals (can also be empty).

Two Weeks Ago: CFGs

A **context-free** grammar (CFG) G consists of:

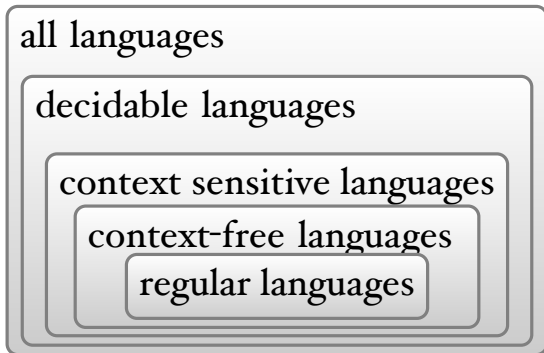
- a finite set of nonterminal symbols (upper case)
- a finite terminal symbols or tokens (lower case)
- a start symbol (which must be a nonterminal)
- a set of rules

$$A \rightarrow \text{rhs}_1 | \text{rhs}_2 | \dots$$

where **rhs** are sequences involving terminals and nonterminals (can also be empty).

Hierarchy of Languages

Recall that languages are sets of strings.



Arithmetic Expressions

A grammar for arithmetic expressions and numbers:

$$\begin{aligned} E &\rightarrow E \cdot + \cdot E \mid E \cdot * \cdot E \mid (\cdot E \cdot) \mid N \\ N &\rightarrow N \cdot N \mid 0 \mid 1 \mid \dots \mid 9 \end{aligned}$$

Unfortunately it is left-recursive (and ambiguous).

A problem for **recursive descent parsers** (e.g. parser combinators).

Arithmetic Expressions

A grammar for arithmetic expressions and numbers:

$$\begin{aligned} E &\rightarrow E \cdot + \cdot E \mid E \cdot * \cdot E \mid (\cdot E \cdot) \mid N \\ N &\rightarrow N \cdot N \mid 0 \mid 1 \mid \dots \mid 9 \end{aligned}$$

Unfortunately it is left-recursive (and ambiguous).

A problem for **recursive descent parsers** (e.g. parser combinators).

Numbers

$$N \rightarrow N \cdot N \mid 0 \mid 1 \mid \dots \mid 9$$

A non-left-recursive, non-ambiguous grammar for numbers:

$$N \rightarrow 0 \cdot N \mid 1 \cdot N \mid \dots \mid 0 \mid 1 \mid \dots \mid 9$$

Operator Precedences

To disambiguate

$$E \rightarrow E \cdot + \cdot E \mid E \cdot * \cdot E \mid (\cdot E \cdot) \mid N$$

Decide on how many precedence levels, say
highest for $()$, medium for $*$, lowest for $+$

$$\begin{aligned} E_{low} &\rightarrow E_{med} \cdot + \cdot E_{low} \mid E_{med} \\ E_{med} &\rightarrow E_{hi} \cdot * \cdot E_{med} \mid E_{hi} \\ E_{hi} &\rightarrow (\cdot E_{low} \cdot) \mid N \end{aligned}$$

Operator Precedences

To disambiguate

$$E \rightarrow E \cdot + \cdot E \mid E \cdot * \cdot E \mid (\cdot E \cdot) \mid N$$

Decide on how many precedence levels, say
highest for $()$, medium for $*$, lowest for $+$

$$\begin{aligned} E_{low} &\rightarrow E_{med} \cdot + \cdot E_{low} \mid E_{med} \\ E_{med} &\rightarrow E_{hi} \cdot * \cdot E_{med} \mid E_{hi} \\ E_{hi} &\rightarrow (\cdot E_{low} \cdot) \mid N \end{aligned}$$

What happens with $1 + 3 + 4$?

Removing Left-Recursion

The rule for numbers is directly left-recursive:

$$N \rightarrow N \cdot N \mid 0 \mid 1 \quad (\dots)$$

Translate

$$\begin{array}{l} N \rightarrow N \cdot \alpha \\ \quad \mid \beta \end{array} \quad \Rightarrow \quad \begin{array}{l} N \rightarrow \beta \cdot N' \\ N' \rightarrow \alpha \cdot N' \\ \quad \mid \epsilon \end{array}$$

Removing Left-Recursion

The rule for numbers is directly left-recursive:

$$N \rightarrow N \cdot N \mid 0 \mid 1 \quad (\dots)$$

Translate

$$\begin{array}{l} N \rightarrow N \cdot \alpha \\ \quad \mid \beta \end{array} \quad \Rightarrow \quad \begin{array}{l} N \rightarrow \beta \cdot N' \\ N' \rightarrow \alpha \cdot N' \\ \quad \mid \epsilon \end{array}$$

Which means

$$\begin{array}{l} N \rightarrow 0 \cdot N' \mid 1 \cdot N' \\ N' \rightarrow N \cdot N' \mid \epsilon \end{array}$$

Chomsky Normal Form

All rules must be of the form

$$A \rightarrow a$$

or

$$A \rightarrow B \cdot C$$

No rule can contain ϵ .

ϵ -Removal

- 1 If $A \rightarrow \alpha \cdot B \cdot \beta$ and $B \rightarrow \epsilon$ are in the grammar, then add $A \rightarrow \alpha \cdot \beta$ (iterate if necessary).
- 2 Throw out all $B \rightarrow \epsilon$.

$$N \rightarrow 0 \cdot N' \mid 1 \cdot N'$$

$$N' \rightarrow N \cdot N' \mid \epsilon$$

$$N \rightarrow 0 \cdot N' \mid 1 \cdot N' \mid 0 \mid 1$$

$$N' \rightarrow N \cdot N' \mid N \mid \epsilon$$

$$N \rightarrow 0 \cdot N' \mid 1 \cdot N' \mid 0 \mid 1$$

$$N' \rightarrow N \cdot N' \mid N$$

ϵ -Removal

- 1 If $A \rightarrow \alpha \cdot B \cdot \beta$ and $B \rightarrow \epsilon$ are in the grammar, then add $A \rightarrow \alpha \cdot \beta$ (iterate if necessary).
- 2 Throw out all $B \rightarrow \epsilon$.

$$N \rightarrow 0 \cdot N' \mid 1 \cdot N'$$

$$N' \rightarrow N \cdot N' \mid \epsilon$$

$$N \rightarrow 0 \cdot N' \mid 1 \cdot N' \mid 0 \mid 1$$

$$N' \rightarrow N \cdot N' \mid N \mid \epsilon$$

$$N \rightarrow 0 \cdot N' \mid 1 \cdot N' \mid 0 \mid 1$$

$$N' \rightarrow N \cdot N' \mid N$$

$$N \rightarrow 0 \cdot N \mid 1 \cdot N \mid 0 \mid 1$$

CYK Algorithm

If grammar is in Chomsky normalform ...

$S \rightarrow N \cdot P$

$P \rightarrow V \cdot N$

$N \rightarrow N \cdot N$

$N \rightarrow \text{students} \mid \text{Jeff} \mid \text{geometry} \mid \text{trains}$

$V \rightarrow \text{trains}$

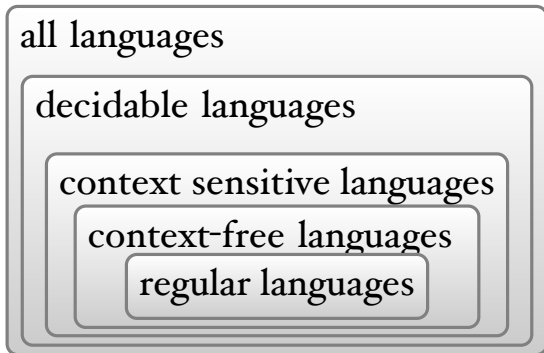
Jeff trains geometry students

CYK Algorithm

- fastest possible algorithm for recognition problem
- runtime is $O(n^3)$
- grammars need to be transferred into CNF

Hierarchy of Languages

Recall that languages are sets of strings.



Context Sensitive Grms

$$S \Rightarrow bSAA \mid \epsilon$$
$$A \Rightarrow a$$
$$bA \Rightarrow Ab$$

Context Sensitive Grms

$$S \Rightarrow bSAA \mid \epsilon$$

$$A \Rightarrow a$$

$$bA \Rightarrow Ab$$

$$S \Rightarrow \dots \Rightarrow^? \text{ "ababaa"}$$

Parse Trees

Stmt → skip
| *Id* := *AExp*
| if *BExp* then *Block* else *Block*
| while *BExp* do *Block*
| read *Id*
| write *Id*
| write *String*

Stmts → *Stmt* ; *Stmts*
| *Stmt*

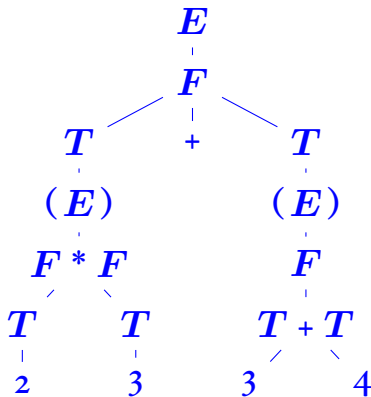
Block → { *Stmts* }
| *Stmt*

AExp → ...

Parse Trees

$$E \rightarrow F \mid F \cdot * \cdot F$$
$$F \rightarrow T \mid T \cdot + \cdot T \mid T \cdot - \cdot T$$
$$T \rightarrow \text{num_token} \mid (.E.)$$

$(2*3)+(3+4)$



Ambiguous Grammars

A CFG is **ambiguous** if there is a string that has at least two parse trees.

$$E \rightarrow \textit{num_token}$$

$$E \rightarrow E \cdot + \cdot E$$

$$E \rightarrow E \cdot - \cdot E$$

$$E \rightarrow E \cdot * \cdot E$$

$$E \rightarrow (\cdot E \cdot)$$

1 + 2 * 3 + 4

Dangling Else

Another ambiguous grammar:

$$\begin{array}{l} E \rightarrow \text{if } E \text{ then } E \\ \quad | \text{if } E \text{ then } E \text{ else } E \\ \quad | \text{id} \end{array}$$

if a then if x then y else c

A CFG Derivation

- 1 Begin with a string with only the start symbol S
- 2 Replace any non-terminal X in the string by the right-hand side of some production $X \rightarrow \text{rhs}$
- 3 Repeat 2 until there are no non-terminals

$S \rightarrow \dots \rightarrow \dots \rightarrow \dots \rightarrow \dots$