



cse 311: foundations of computing

Spring 2015

Lecture 20: Regular expressions and context-free grammars



languages: sets of strings

Sets of strings that satisfy special properties are called **languages**.

Examples:

- English sentences
- Syntactically correct Java/C/C++ programs
- Σ^* = All strings over alphabet Σ
- Palindromes over Σ
- Binary strings that don't have a 0 after a 1
- Legal variable names, keywords in Java/C/C++
- Binary strings with an equal # of 0's and 1's

regular expressions

Regular expressions over Σ

- Basis:
 - \emptyset, ϵ are regular expressions
 - a is a regular expression for any $a \in \Sigma$
- Recursive step:
 - If **A** and **B** are regular expressions then so are:
 - $(A \cup B)$
 - (AB)
 - A^*

each regular expression is a "pattern"

 ϵ matches the **empty string** a matches the one character string a $(A \cup B)$ matches all strings that either **A** matches or **B** matches (or both) (AB) matches all strings that have a first part that **A** matches followed by a second part that **B** matches A^* matches all strings that have any number of strings (even 0) that **A** matches, one after another

examples

- 001^*
- 0^*1^*
- $(0 \cup 1)0(0 \cup 1)0$
- $(0^*1^*)^*$
- $(0 \cup 1)^*0110(0 \cup 1)^*$
- $(00 \cup 11)^*(01010 \cup 10001)(0 \cup 1)^*$

regular expressions in practice

- Used to define the "tokens": e.g., legal variable names, keywords in programming languages and compilers
- Used in **grep**, a program that does pattern matching searches in UNIX/LINUX
- Pattern matching using regular expressions is an essential feature of PHP
- We can use regular expressions in programs to process strings!

example

Example: $S \rightarrow OS0 \mid 1S1 \mid 0 \mid 1 \mid \varepsilon$

Example: $S \rightarrow OS \mid S1 \mid \varepsilon$

example

Grammar for $\{0^n 1^n : n \geq 0\}$

(all strings with same # of 0's and 1's with all 0's before 1's)

Example: $S \rightarrow (S) \mid SS \mid \varepsilon$

simple arithmetic expressions

$E \rightarrow E+E \mid E*(E) \mid x \mid y \mid z \mid 0 \mid 1 \mid 2 \mid 3 \mid 4$
 $\mid 5 \mid 6 \mid 7 \mid 8 \mid 9$

Generate $(2*x) + y$

Generate $x+y*z$ in two fundamentally different ways

parse trees

Suppose that grammar G generates a string x

A **parse tree** of x for G has

- Root labeled S (start symbol of G)
- The children of any node labeled A are labeled by symbols of w left-to-right for some rule $A \rightarrow w$
- The symbols of x label the leaves ordered left-to-right

$S \rightarrow OS0 \mid 1S1 \mid 0 \mid 1 \mid \varepsilon$

Parse tree of **01110**:



CFGs and recursively-defined sets of strings

- A CFG with the start symbol S as its only variable recursively defines the set of strings of terminals that S can generate
- A CFG with more than one variable is a simultaneous recursive definition of the sets of strings generated by *each* of its variables
 - Sometimes necessary to use more than one

building precedence in simple arithmetic expressions

- E – expression (start symbol)
- T – term F – factor I – identifier N – number
 - $E \rightarrow T \mid E+T$
 - $T \rightarrow F \mid F*T$
 - $F \rightarrow (E) \mid I \mid N$
 - $I \rightarrow x \mid y \mid z$
 - $N \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$

Backus-Naur form (same as CFG)

BNF (Backus-Naur Form) grammars

- Originally used to define programming languages
- Variables denoted by long names in angle brackets, e.g.
 <identifier>, <if-then-else-statement>,
 <assignment-statement>, <condition>
 ::= used instead of →

BNF for C

```
statement:
  ((identifier | "case" constant-expression | "default") ":"*)
  (expression? ";" |
   block |
   "if" "(" expression ")" statement |
   "if" "(" expression ")" statement "else" statement |
   "switch" "(" expression ")" statement |
   "while" "(" expression ")" statement |
   "do" statement "while" "(" expression ")" ";" |
   "for" "(" expression? ";" expression? ";" expression? ")" statement |
   "goto" identifier ";" |
   "continue" ";" |
   "break" ";" |
   "return" expression? ";"
  )

block: "(" declaration* statement* ")"

expression:
  assignment-expression

assignment-expression: (
  unary-expression (
    "=" | "&=" | "/=" | "%=" | "+=" | "-=" | "<<=" | ">>=" | "=" |
    "&=" | "|="
  )
)* conditional-expression

conditional-expression:
  logical-OR-expression ( "?" expression ":" conditional-expression )?
```

parse trees

Back to middle school:

<sentence>::=<noun phrase><verb phrase>
 <noun phrase>::=<article><adjective><noun>
 <verb phrase>::=<verb><adverb><verb><object>
 <object>::=<noun phrase>

Parse:

The yellow duck squeaked loudly
 The red truck hit a parked car