CSE 311: Foundations of Computing

Topic 7: Languages



Theoretical Computer Science

- An alphabet Σ is any finite set of characters
- The set Σ^* of strings over the alphabet Σ
 - example: {0,1}* is the set of binary strings
 0, 1, 00, 01, 10, 11, 000, 001, ... and ""
- Σ^* is defined recursively by
 - Basis: $\varepsilon \in \Sigma^*$ (ε is the empty string, i.e., "")
 - **Recursive:** if $w \in \Sigma^*$, $a \in \Sigma$, then $wa \in \Sigma^*$

- Subsets of strings are called languages
- Examples:
 - $-\Sigma^* = \text{All strings over alphabet } \Sigma$
 - palindromes over $\boldsymbol{\Sigma}$
 - binary strings with an equal # of 0's and 1's
 - syntactically correct Java/C/C++ programs
 - valid English sentences
 - correct solutions to coding problems:

 $S = \{x # y | y \text{ is Java code that does what } x \text{ says}\}$

- Look at different ways of defining languages
- See which are more expressive than others
 - i.e., which can define more languages
- Later: connect ways of defining languages to different types of (restricted) computers
 - computers capable of recognizing those languages
 i.e., distinguishing strings in the language from not
- Consequence: computers that recognize more expressive languages are more powerful

Palindromes are strings that are the same when read backwards and forwards

Basis:

 ε is a palindrome any $a \in \Sigma$ is a palindrome

Recursive step:

If p is a palindrome, then apa is a palindrome for every $a \in \Sigma$

Regular expressions over $\boldsymbol{\Sigma}$

• Basis:

ε is a regular expression (could also include ∅) α is a regular expression for any α ∈ Σ

• Recursive step:

If **A** and **B** are regular expressions, then so are:

A ∪ B AB A*

- ε matches only the empty string
- *a* matches only 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 ($\varepsilon \cup A \cup AA \cup AA \cup ...$)

Definition of the *language* matched by a regular expression The language defined by a regular expression:

$$L(\varepsilon) = \{\varepsilon\}$$

$$L(a) = \{a\}$$

$$L(A \cup B) = L(A) \cup L(B)$$

$$L(AB) = \{y \bullet z : y \in L(A), z \in L(B)\}$$

$$L(A^*) = \bigcup_{n=0}^{\infty} L(A^n)$$

$$A^n \text{ defined recursively by}$$

$$A^0 = \emptyset$$

$$A^{n+1} = A^n A$$

001*

0*1*

001*

 $\{00, 001, 0011, 00111, ...\}$

0*1*

Any number of 0's followed by any number of 1's

 $(\mathbf{0} \cup \mathbf{1}) \, \mathbf{0} \, (\mathbf{0} \cup \mathbf{1}) \, \mathbf{0}$



 $(\mathbf{0} \cup \mathbf{1}) \, \mathbf{0} \, (\mathbf{0} \cup \mathbf{1}) \, \mathbf{0}$

 $\{0000, 0010, 1000, 1010\}$

(0*1*)*

All binary strings

• All binary strings that contain 0110

```
(0 \cup 1)* 0110 (0 \cup 1)*
```

• All binary strings that begin with a string of doubled characters (00 or 11) followed by 01010 or 10001

 $(00 \cup 11)*(01010 \cup 10001)(0 \cup 1)*$

• All binary strings that have an even # of 1's

e.g., 0*(10*10*)*

• All binary strings that *don't* contain 101

e.g., 0*(1 U 1000*)*(ε U 10)

at least two 0s between 1s

Finite languages vs Regular Expressions

• All finite languages have a regular expression. (a language is finite if its elements can be put into a list)

Why?

• Given a list of strings $s_1, s_2, ..., s_n$

Construct the regular expression

 $s_1 U s_2 U \dots U s_n$

(Could make this formal by induction on n)

Finite languages vs Regular Expressions

• Every regular expression that does not use * generates a finite language.

Why?

• Prove by structural induction on the syntax of regular expressions!

Let A be a regular expression that does not use *. Then L(A) is finite.

Proof: We proceed by structural induction on A.

Case \epsilon: $L(\epsilon) = \{\epsilon\}$, which is finite

Case a: $L(a) = \{a\}$, which is finite

Case A \cup B: L(A \cup B) = L(A) \cup L(B) By the IH, each is finite, so their union is finite. Let A be a regular expression that does not use *. Then L(A) is finite.

Proof: We proceed by structural induction on A. Case AB: $L(AB) = \{y \bullet z : y \in L(A), z \in L(B)\}$ By the IH, L(A) and L(B) are finite.

Every element of L(AB) is covered by a pair (y, z) where $y \in L(A)$ and $z \in L(B)$, so L(AB) is finite.

(No case for A*!)

Finite languages vs Regular Expressions

Key takeaways:

- Regular expressions can represent all finite languages
- To prove a language is "regular", just give the regular expression that describes it.
- Regular expressions are more powerful than finite languages (e.g., 0* is an infinite language)
- To prove something about *all* regular expressions, use structural induction on the syntax.

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!

Regular Expressions in Java

- Pattern p = Pattern.compile("a*b");
- Matcher m = p.matcher("aaaaab");
- boolean b = m.matches();
 - [01] a 0 or a 1 ^ start of string \$ end of string
 - [0-9] any single digit $\$. period $\$, comma $\$ minus . any single character
 - ab a followed by b (AB)
 - (a|b) a or b $(A \cup B)$
 - a? zero or one of a $(\mathbf{A} \cup \boldsymbol{\varepsilon})$
 - a* zero or more of a A*
 - a+ one or more of a **AA***

e.g. ^[\-+]?[0-9]*(\.|\,)?[0-9]+\$
 General form of decimal number e.g. 9.12 or -9,8 (Europe)

Limitations of Regular Expressions

- Not all languages can be specified by regular expressions
- Even some easy things like
 - Palindromes
 - Strings with equal number of 0's and 1's
- But also more complicated structures in programming languages
 - Matched parentheses
 - Properly formed arithmetic expressions
 - etc.

Example: $S \rightarrow 0S0 \mid 1S1 \mid 0 \mid 1 \mid \epsilon$

Example: $S \rightarrow 0S0 | 1S1 | 0 | 1 | \epsilon$

How does this grammar generate 0110?

Example: $S \rightarrow 0S0 \mid 1S1 \mid 0 \mid 1 \mid \epsilon$

How does this grammar generate 0110?

$\textbf{S} \rightarrow 0\textbf{S}0 \rightarrow 01\textbf{S}10 \rightarrow 01\epsilon10 = 0110$

Example: $\mathbf{S} \rightarrow \mathbf{0S0} \mid \mathbf{1S1} \mid \mathbf{0} \mid \mathbf{1} \mid \mathbf{\epsilon}$

How to describe all strings generated?

The set of all binary palindromes

Recursively-Defined Set

Basis:

ε is a palindrome any $a \in \{0, 1\}$ is a palindrome

Recursive step:

If *p* is a palindrome, then apa is a palindrome for every $a \in \{0, 1\}$

$Grammar \qquad S \rightarrow 0S0 \mid 1S1 \mid 0 \mid 1 \mid \epsilon$

Example: $S \rightarrow A \mid B$ $A \rightarrow 0A \mid \varepsilon$ $B \rightarrow 1B \mid \varepsilon$

How does this grammar generate 000?

Example: $S \rightarrow A \mid B$ $A \rightarrow 0A \mid \varepsilon$ $B \rightarrow 1B \mid \varepsilon$

How does this grammar generate 000?

 $\textbf{S} \rightarrow \textbf{A} \rightarrow 0\textbf{A} \rightarrow 00\textbf{A} \rightarrow 000\textbf{A} \rightarrow 000\epsilon = 000$

Example: $S \rightarrow A \mid B$ $A \rightarrow 0A \mid \varepsilon$ $B \rightarrow 1B \mid \varepsilon$

How to describe all strings generated?

strings of all 0s or all 1s

(all 0s) ∪ (all 1s)

- A Context-Free Grammar (CFG) is given by a finite set of substitution rules involving
 - A finite set V of variables that can be replaced
 - Alphabet Σ of *terminal symbols* that can't be replaced
 - One variable, usually **S**, is called the *start symbol*
- The substitution rules involving a variable **A**, written as $\begin{array}{c|c} \mathbf{A} \to w_1 & w_2 & \cdots & w_k \\ \hline w_1 & w_2 & \cdots & w_k \\ \hline w_1 & w_2 & w_1 & w_2 \\ \hline w_1 & w_2 & \cdots & w_k \end{array}$

- that is $w_i \in (\mathbf{V} \cup \Sigma)^*$

- Begin with start symbol **S**
- If there is some variable **A** in the current string you can replace it by one of the w's in the rules for **A**

$$- \mathbf{A} \rightarrow \mathbf{w}_1 \mid \mathbf{w}_2 \mid \cdots \mid \mathbf{w}_k$$

- Write this as $xAy \Rightarrow xwy$
- Repeat until no variables left
- The set of strings the CFG describes are all strings, containing no variables, that can be *generated* in this manner (after a finite number of steps)

Example: $S \rightarrow 0S \mid S1 \mid \epsilon$

Example: $S \rightarrow 0S | S1 | \epsilon$

0*1*

(i.e., matching 0*1* but with same number of 0's and 1's)

(i.e., matching 0*1* but with same number of 0's and 1's)

$\textbf{S} \rightarrow \textbf{OS1} ~|~ \epsilon$

(i.e., matching 0*1* but with same number of 0's and 1's)

$\textbf{S} \rightarrow \textbf{OS1} ~|~ \epsilon$

Grammar for $\{0^n 1^{2n} : n \ge 0\}$

(i.e., matching 0*1* but with same number of 0's and 1's)

$\textbf{S} \rightarrow \textbf{OS1} ~|~ \epsilon$

Grammar for $\{0^n 1^{2n} : n \ge 0\}$

$S \rightarrow 0S11 \mid \epsilon$

(i.e., matching 0*1* but with same number of 0's and 1's)

$\textbf{S} \rightarrow \textbf{OS1} ~|~ \epsilon$

Grammar for $\{0^n 1^{n+1} 0 : n \ge 0\}$

(i.e., matching 0*1* but with same number of 0's and 1's)

$\textbf{S} \rightarrow \textbf{OS1} ~|~ \epsilon$

Grammar for $\{0^n 1^{n+1} 0 : n \ge 0\}$

 $S \rightarrow A 10$ $A \rightarrow 0A1 | \epsilon$

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

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

The set of all strings of matched parentheses

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

The set of all strings of matched parentheses

Suppose S generates x. Define f(k) to be number of "("s – ")"s in first k characters of x



Example Context-Free Grammars

Three possibilities for f(k) for $k \in \{1, ..., n-1\}$

• f(k) > 0 for all such k $S \rightarrow (S)$



• f(k) = 0 for some such k

 $S \rightarrow SS$



Binary strings with equal numbers of 0s and 1s (not just 0ⁿ1ⁿ, also 0101, 0110, etc.)

 $\textbf{S} \rightarrow \textbf{SS}$ | 0S1 | 1S0 | ϵ

Binary strings with equal numbers of 0s and 1s (not just 0ⁿ1ⁿ, also 0101, 0110, etc.)

 $\textbf{S} \rightarrow \textbf{SS}$ | 0S1 | 1S0 | ϵ

Suppose S generates x. Define f(k) to be #0s – #1s in the first k characters of x.



Binary strings with equal numbers of 0s and 1s (not just 0ⁿ1ⁿ, also 0101, 0110, etc.)

$\textbf{S} \rightarrow \textbf{SS} \mid \textbf{0S1} \mid \textbf{1S0} \mid \epsilon$

Suppose S generates x. Define f(k) to be #0s – #1s in the first k characters of x.

If k-th character is 0, then f(k) = f(k-1) + 1If k-th character is 1, then f(k) = f(k-1) - 1 Let $x \in (0 \cup 1)^*$. Define $f_x(k)$ to be the number Os minus the number of 1s in the k characters of x.



f(k) = 0 when first k characters have #0s = #1s

- starts out at 0f(0) = 0- ends at 0f(n) = 0

Three possibilities for f(k) for $k \in \{1, ..., n-1\}$

- f(k) > 0 for all such k**S** \rightarrow **0S1**
- f(k) < 0 for all such k

 $\mathbf{S}
ightarrow \mathbf{1S0}$

• f(k) = 0 for some such k

 $S \rightarrow SS$





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

Generate (2 + x) * y

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

Generate (2 + x) * y

 $\mathsf{E} \Rightarrow \mathsf{E}^*\mathsf{E} \Rightarrow (\mathsf{E})^*\mathsf{E} \Rightarrow (\mathsf{E}+\mathsf{E})^*\mathsf{E} \Rightarrow (\mathsf{2}+\mathsf{E})^*\mathsf{E} \Rightarrow (\mathsf{2}+x)^*\mathsf{E} \Rightarrow (\mathsf{2}+x)^*\mathsf{y}$

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

 $\mathbf{S} \rightarrow \mathbf{0S0} \mid \mathbf{1S1} \mid \mathbf{0} \mid \mathbf{1} \mid \mathbf{\epsilon}$



Parse tree of 01110

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

Generate x+y*z in two ways that give two *different* parse trees

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

Generate x+y*z in ways that give two *different* parse trees



Structural induction is the tool used to prove many more interesting theorems

- General associativity follows from our one rule
 - likewise for generalized De Morgan's laws
- Okay to substitute y for x everywhere in a modular equation when we know that $x \equiv_m y$
- The "Meta Theorem" on set operators

These are proven by induction on parse trees

parse trees are recursively defined

Theorem: For any set of strings (language) *A* described by a regular expression, there is a CFG that recognizes *A*.

Proof idea:

P(A) is "A is recognized by some CFG"

Structural induction based on the recursive definition of regular expressions...

• Basis:

- $-\epsilon$ is a regular expression
- **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***

CFGs are more general than **REs**

• CFG to match RE **E**

 $\textbf{S} \rightarrow \epsilon$

• CFG to match RE **a** (for any $a \in \Sigma$)

 $S \rightarrow a$

CFGs are more general than **REs**

Suppose CFG with start symbol **S**₁ matches RE **A** CFG with start symbol **S**₂ matches RE **B**

- CFG to match RE $\mathbf{A} \cup \mathbf{B}$
 - $S \rightarrow S_1 \mid S_2$ + rules from original CFGs
- CFG to match RE **AB**

 $\mathbf{S} \rightarrow \mathbf{S}_1 \mathbf{S}_2$ + rules from original CFGs

CFGs are more general than **REs**

Suppose CFG with start symbol S_1 matches RE A

• CFG to match RE A^* (= $\varepsilon \cup A \cup AA \cup AAA \cup ...$)

 $S \rightarrow S_1 S \mid \epsilon$ + rules from CFG with S_1