

Computational Models — Lecture 4¹

Handout Mode

Ronitt Rubinfeld and Iftach Haitner.

Tel Aviv University.

March 30/April 13, 2015

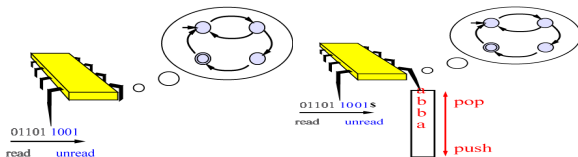
¹Based on frames by Benny Chor, Tel Aviv University, modifying frames by Maurice Herlihy, Brown University.

Talk Outline

- ▶ **Context Free** Grammars/Languages (CFG/CFL)
- ▶ (begin) Algorithmic issues for CFL

Next two weeks:

- ▶ Finish algorithmic issues for CFL
- ▶ **Pumping Lemma** for context free languages
- ▶ Push Down Automata (**PDA**)
- ▶ **Equivalence** of CFGs and PDAs



- ▶ Sipser's book, 2.1, 2.2 & 2.3

Short Overview of the Course

So far we saw

- ▶ finite automata,
- ▶ regular languages,
- ▶ regular expressions,
- ▶ Myhill-Nerode theorem
- ▶ pumping lemma for **regular languages**.

We now introduce stronger machines and languages with more expressive power:

- ▶ pushdown automata,
- ▶ context-free languages,
- ▶ context-free grammars,
- ▶ pumping lemma for context-free languages.

Context Free Grammars (CFG)

An example of a context free grammar, G_1 :

- ▶ $A \rightarrow 0A1$
- ▶ $A \rightarrow B$
- ▶ $B \rightarrow \#$

Terminology:

- ▶ Each line is a **substitution rule** or **production**.
- ▶ Each rule has the form: **symbol** \rightarrow **string**.
The left-hand symbol is a **variable** (usually upper-case).
- ▶ A **string** consists of **variables** and **terminals**.
- ▶ One variable is the **start variable** (lhs of top rule). In this case, it is A .

Rules for Generating Strings

- ▶ Write down the start variable.
- ▶ Pick a variable written down in current string and a derivation that starts with that variable.
- ▶ Replace that variable with right-hand side of that derivation.
- ▶ Repeat until no variables remain.
- ▶ Return final string (concatenation of terminals).

Process is inherently **non deterministic**.

Example

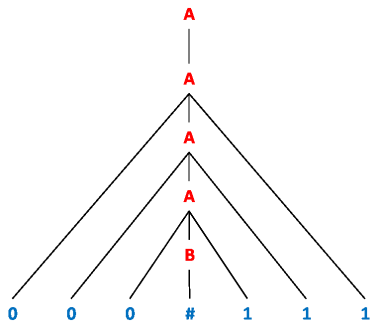
Grammar G_1 :

- ▶ $A \rightarrow 0A1$
- ▶ $A \rightarrow B$
- ▶ $B \rightarrow \#$

Derivation with G_1 :

A \rightarrow 0A1
 \rightarrow 00A11
 \rightarrow 000A111
 \rightarrow 000B111
 \rightarrow **000#111**

Parsing Trees



Indifferent to derivation **order**.

Question 1

What strings can be generated in this way from the grammar G_1 ?

Answer: Exactly those of the form $0^n \# 1^n$ ($n \geq 0$).

Context-Free Languages (CFL)

The language generated in this way is called the **language of the grammar**.

For example, $\mathcal{L}(G_1) = \{0^n \# 1^n : n \geq 0\}$.

Any language generated by a context-free grammar is called a **context-free language**.

A Useful Abbreviation

Rules with same variable on left hand side

$$A \rightarrow 0A1$$

$$A \rightarrow B$$

are written as:

$$A \rightarrow 0A1 \mid B$$

English-like Sentences

A grammar G_2 to describe a few English sentences:

$\langle \text{SENTENCE} \rangle \rightarrow \mathcal{NP} \langle \text{VERB} \rangle$
 $\mathcal{NP} \rightarrow \langle \text{ARTICLE} \rangle \langle \text{NOUN} \rangle$
 $\langle \text{NOUN} \rangle \rightarrow \text{boy} \mid \text{girl} \mid \text{flower}$
 $\langle \text{ARTICLE} \rangle \rightarrow \text{a} \mid \text{the}$
 $\langle \text{VERB} \rangle \rightarrow \text{touches} \mid \text{likes} \mid \text{sees}$

Deriving English-like Sentences

A specific derivation in G_2 :

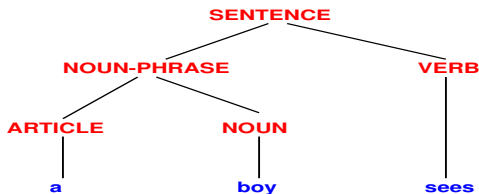
< SENTENCE > \rightarrow \mathcal{NP} < VERB >
 \rightarrow < ARTICLE >< NOUN >< VERB >
 \rightarrow a < NOUN >< VERB >
 \rightarrow a boy < VERB >
 \rightarrow a boy sees

More strings generated by G_2 :

a flower sees
the girl touches

Derivation and Parse Tree

- $\langle \text{SENTENCE} \rangle \rightarrow \mathcal{NP} \langle \text{VERB} \rangle$
- $\langle \text{ARTICLE} \rangle \langle \text{NOUN} \rangle \langle \text{VERB} \rangle$
- **a** $\langle \text{NOUN} \rangle \langle \text{VERB} \rangle$
- **a boy** $\langle \text{VERB} \rangle$
- **a boy sees**



Formal Definition

A **context-free grammar** is a 4-tuple (V, Σ, R, S) , where

- ▶ V is a finite set of **variables**
- ▶ Σ is a finite set of **terminals**
- ▶ R is a finite set of **rules**: each rule is a variable and a finite string of variables and terminals.
- ▶ S is the **start symbol**.
- ▶ If u and v are strings of **variables** and **terminals**, and $A \rightarrow w$ is a **rule** of the grammar, then uAv **yields** uwv , written $uAv \rightarrow uwv$.
- ▶ $u \xrightarrow{*} v$ if $u = v$, or $u \rightarrow u_1 \rightarrow \dots \rightarrow u_k \rightarrow v$ for some sequence u_1, u_2, \dots, u_k

Definition 2

The **language of the grammar** G , denoted $\mathcal{L}(G)$, is $\{w \in \Sigma^* : S \xrightarrow{*} w\}$

where $\xrightarrow{*}$ is determined by G .

Example 1

$G_3 = (\{S\}, \{a, b\}, R, S)$.

R (Rules): $S \rightarrow aSb \mid SS \mid \epsilon$

Some words in the language: $aabb$, $aababb$.

Question 3

What **is** this language?

Hint: Think of parentheses: i.e., a is "(" and b is ")". $(())$, $((()))$

Using larger alphabet (i.e., more terminals), $([]())$, represent well formed programs with many kinds of nested loops, "if then/else" statements.

Example 2

$$G_4 = (\{S\}, \{a, b\}, R, S).$$

R (Rules): $S \rightarrow aSa \mid bSb \mid \epsilon$

Some words in the language: *abba, aabaabaa*.

Question 4

What **is** this language?

$\mathcal{L}(G_4) = \{ww^R : w \in \{a, b\}^*\}$ (almost but not quite the set of palindromes)

Proving $\mathcal{L}(G_4) = \{ww^R : w \in \{a, b\}^*\}$

What do we need to show?

1. If $z = ww^R$ then z is generated by G_4 .
2. If z generated by G_4 then $z = ww^R$.

Part 1:

Proof by induction on the length of z .

Base: $|z| = 0$, then $S \rightarrow \varepsilon$. Hence $\varepsilon \in \mathcal{L}(G_4)$.

Inductive claim:

Let $|z| = 2k$.

Let $z = ww^R$ and $w = \sigma w'$. Then $z = \sigma w'(w')^R \sigma$.

Consider the derivation $S \rightarrow \sigma S \sigma$.

By the induction hypothesis $S \xrightarrow{*} w'(w')^R$, since $|w'(w')^R| = 2k - 2 < |z|$.

Therefore, G_5 generates $S \rightarrow \sigma S \sigma \xrightarrow{*} \sigma w'(w')^R \sigma = z$.

Hence $z \in \mathcal{L}(G_4)$.

Proving $\mathcal{L}(G_4) = \{ww^R : w \in \{a, b\}^*\}$

Part 2:

Assume that $S \xrightarrow{*} z$.

We prove by induction on the number of derivations that $z = ww^R$.

Base: If we have a single derivation, then the only possible derivation is $S \rightarrow \varepsilon$, and we are done.

Inductive claim:

Assume $S \xrightarrow{*} z$ in k derivations.

Assume we have $S \rightarrow \sigma S \sigma$ as the first derivation.

Then by the inductive hypothesis, $S \rightarrow w'(w')^R$.

Therefore, $z = \sigma w'(w')^R \sigma$.

Example 3

$$G_5 = (\{S, A, B\}, \{a, b\}, R, S)$$

R (Rules):

$$S \rightarrow aB \mid bA \mid \varepsilon$$

$$A \rightarrow a \mid aS \mid bAA$$

$$B \rightarrow b \mid bS \mid aBB$$

Some words in the language: *aababb*, *baabba*.

Question 5

What is this language?

$$\mathcal{L}(G_5) = \{w \in \{a, b\}^* : \#_a(w) = \#_b(w)\}$$

$\#_x(w)$ — number of occurrences of x in w

Proving $\mathcal{L}(G_5) = \{w \in \{a, b\}^* : \#_a(w) = \#_b(w)\}$

What do we need to show?

1. If w generated by G_5 then $\#_a(w) = \#_b(w)$.
2. If $\#_a(w) = \#_b(w)$ then w is generated by G_5 .

Claim 6

If $S \xrightarrow{*} w$, then $\#_a(w) + \#_A(w) = \#_b(w) + \#_B(w)$.

Claim 7

For $w \in \{a, b\}^*$ let $k = k(w) = \#_a(w) - \#_b(w)$. Then:

1. If $k = 0$, then $S \xrightarrow{*} wS$
2. If $k > 0$, then $S \xrightarrow{*} wB^k$
3. If $k < 0$, then $S \xrightarrow{*} wA^{|k|}$

Are we done?! How we prove these claims?

Proving Claim 7

Proof by induction on $|w|$:

- ▶ Basis: $w = \epsilon$ then $k(w) = 0$. Since $S \xrightarrow{*} S$, then $S \xrightarrow{*} \epsilon S$
- ▶ Induction step: for $w \in \{a, b\}^n$ write $w = w'\sigma$ with $|w'| = n - 1$

Suppose:

1. $k' = k(w') = (\#_a(w') - \#_b(w')) > 0$
2. $\sigma = a$

Note that if both are the case, we have $k' = k - 1$

By the induction hypothesis $S \xrightarrow{*} w' B^{k'}$

Since $B \rightarrow aBB$, then

$$S \xrightarrow{*} w' B^{k'} = w' B B^{k'-1} \xrightarrow{*} w' a B B B^{k'-1} = w B^{k'+1} = w B^{k(w)}$$

To complete the proof, need to show for $\sigma = b$ and for $k' \leq 0$

Designing Context-Free Grammars

No recipe in general, but few rules-of-thumb

- ▶ If CFG is the **union** of several CFGs, rename variables (**not terminals**) so they are disjoint, and add new rule $S \rightarrow S_1 \mid S_2 \mid \dots \mid S_j$.
- ▶ For languages (like $\{0^n \# 1^n : n \geq 0\}$), with **linked** substrings, a rule of form $R \rightarrow uRv$ is helpful to force desired relation between substrings.
- ▶ For a regular language, grammar “follows” a DFA for the language (see next frame).

How expressive are CFG's

Are they more expressive or less expressive than regular languages?

CFG for Regular Languages

Given a DFA : $M = (Q, \Sigma, \delta, q_0, F)$

CFG G for $\mathcal{L}(M)$:

1. Let R_0 be the starting variable
2. Add rule $R_i \rightarrow aR_j$, for any $q_i, q_j \in Q$ and $a \in \Sigma$ with $\delta(q_i, a) = q_j$
3. Add rule $R_i \rightarrow \varepsilon$, for any $q_i \in F$

Claim 8

$$\mathcal{L}(G) = \mathcal{L}(M)$$

Proof?

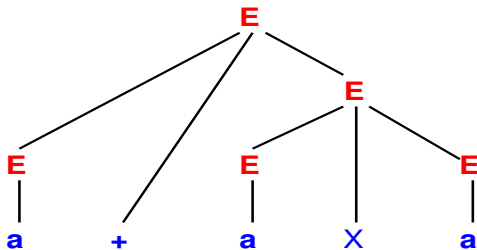
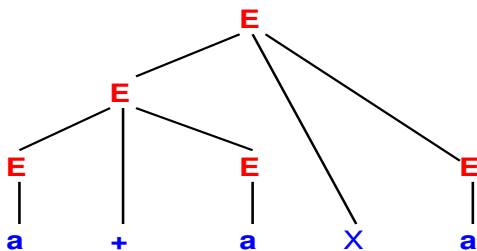
Claim 9

$$R_0 \xrightarrow{*} wR_j \text{ iff } M(w) = q_j.$$

Proof? Next class we'll see alternative proof via "Push-Down Automata"

Ambiguity in CFLs

Grammar $G: E \rightarrow E + E \mid E \times E \mid (E) \mid a$



Arithmetic Example

Consider the grammar $G' = (V, \Sigma, R, E)$, where

▶ $V = \{E, T, F\}$

▶ $\Sigma = \{a, +, \times, (,)\}$

▶ Rules:
$$\begin{array}{l} E \rightarrow E + T \mid T \\ T \rightarrow T \times F \mid F \\ F \rightarrow (E) \mid a \end{array}$$

Claim 10

$$\mathcal{L}(G') = \mathcal{L}(G)$$

Proof? but G' is not ambiguous.

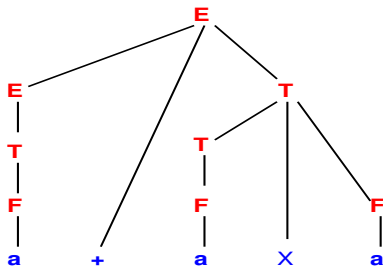
Parsing Tree of G' for $a + a \times a$

G' :

$$E \rightarrow E + T \mid T$$

$$T \rightarrow T \times F \mid F$$

$$F \rightarrow (E) \mid a$$



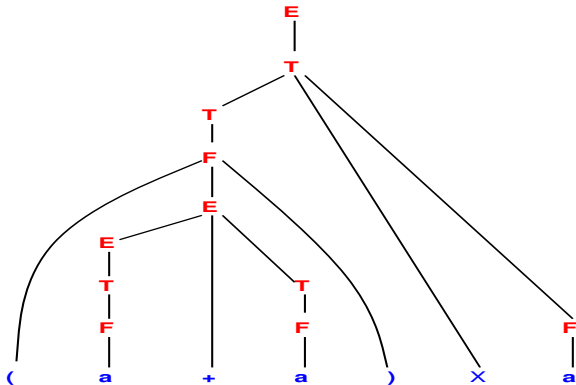
Parsing Tree of G' for $(a + a) \times a$

G' :

$$E \rightarrow E + T \mid T$$

$$T \rightarrow T \times F \mid F$$

$$F \rightarrow (E) \mid a$$



Ambiguity

Definition 11

A string w is derived **ambiguously** from grammar G , if w has two or more **different** parse trees that generate it from G . A CFG is **ambiguous**, if it ambiguously derives a string.

- ▶ Ambiguity is usually not only a syntactic notion but also **semantic**, implying multiple meanings for the same string. Think of $a + a \times a$ from last grammar.
- ▶ It is **sometime** possible to **eliminate** ambiguity by finding a different context free grammar generating the same language. This is true for the **arithmetic expressions** grammar.

Some languages are **inherently** ambiguous.

Example: $\{1^i 2^j 3^k : i = j \vee j = k\}$

Proof? Book!

Part I

Checking Membership

Checking Membership in a CFL

Challenge

Given a CFG G and a string w , decide whether $w \in \mathcal{L}(G)$?

Initial Idea: Design an algorithm that tries **all derivations**.

Problem: If G does **not** generate w , we'll never stop.

Possible solution: Use special grammars that are:

- ▶ just as expressive!
- ▶ better for checking membership.