

COL352 Lecture 14

Contents

1	Recap	1
2	Definitions	1
3	Content	1
3.1	Context-Free Languages	1

1 Recap

Discussion in previous class about Myhill-Nerode theorem and starting of context free languages.

2 Definitions

3 Content

3.1 Context-Free Languages

Motivating example:

Example 1

Inductively define L as the smallest language satisfying the following:

1. $\epsilon \in L$.
2. If $x, y \in L$ then $x \cdot y \in L$.
3. If $x \in L$, then $0x1 \in L$.

Note 1

This is something like balanced parenthesized expressions. We needed the smallest language thing because any superset of L also works.

Claim 0.1

$x \in L \iff$ the number of 0s in x = number of 1s in x , and for all y (prefixes of x), the number of 0s in y is at least the number of 1s in y .

Proof. Exercise. □

Claim 0.2

This language is not regular.

Proof. Use the pumping lemma on $0^n 1^n$ or give another proof using the Myhill-Nerode theorem. □

This is something like a grammar.

In what follows, S is the initial non-terminal.

1. $S \rightarrow \epsilon$

2. $S \rightarrow SS$
3. $S \rightarrow 0S1$

0,1 are called terminals, S is a non-terminal, the above three rules are called the production rules for the grammar.

We can make parse trees using this.

Example 2

We want to make a primitive calculator to add and multiply non-negative integers. $\Sigma = \{0, \dots, 9, +, *, (,)\}$

$E \rightarrow E + T | T$

$T \rightarrow T * F | F$

$F \rightarrow (E) | N$

$N \rightarrow 0N | 1N | \dots | 9N | 0 | 1 | \dots | 9$

Here E stands for an expression, T for a term, F for a factor, N for a number.

$X \rightarrow Y | Z$ is shorthand for the two rules $X \rightarrow Y$ and $X \rightarrow Z$.

Definition 1

A grammar (context-free grammar) G is a 4-tuple of the form (N, Σ, R, S) where

1. N = a finite nonempty set of “nonterminals” (a.k.a. “variables”),
2. Σ = a finite nonempty alphabet (set of terminals)
3. $R \subseteq N \times ((N \cup \Sigma)^*)$ = the set of production rules
4. $S \in N$ = initial non-terminal

Note that $(X, w) \in R$ is also written as $X \rightarrow_G w$ or $X \rightarrow w$ if G is clear from context.

Note that this is analogous to regular expressions (i.e., this is not a machine).

We shall make the ‘expansion’ of a non-terminal more formal in the following definition.

Definition 2

Let $G = (N, \Sigma, R, S)$ be a grammar, and $x, y \in (N \cup \Sigma)^*$. We say v produces y and write $x \Rightarrow_G y$ if $\exists X \in N, u, v, w \in (N \cup \Sigma)^*$ such that $x = uXv, y = uvw, X \rightarrow_G w \in R$.

Definition 3

We say “ x derives y ”, and write $x \xRightarrow{*}_G y$ if \exists a finite sequence of strings $z_0, \dots, z_n \in (N \cup \Sigma)^*$ such that $x = z_0, y = z_n$ and $\forall i, z_i \Rightarrow z_{i+1}$.

Note that the derivation relation is the reflexive transitive closure of the production relation.

Definition 4

The language of grammar $G = (N, \Sigma, R, S)$, written as $\mathcal{L}(G)$ is the set of all $x \in \Sigma^*$ such that $S \xRightarrow{*}_G x$.

Example 3

$E \Rightarrow E + T \Rightarrow T + T \Rightarrow T + T * F \Rightarrow \dots \Rightarrow F + F * F \Rightarrow \dots \Rightarrow 1 + 2 * 3$ - the derivation of $1 + 2 * 3$.

Example 4

Using the first example we did, we can do $S \Rightarrow 0S1 \Rightarrow 0SS1 \xRightarrow{*} 00S10S11 \xRightarrow{*} 001011$.

The parse tree’s yield is what you get from reading the parse tree’s leaves in left to right order (parse tree doesn’t need to have leaves as terminals).

Definition 5

Let $G = (N, \Sigma, R, S)$ be a grammar. The set of parse trees of G and the root and yield of each parse tree is recursively defined as follows:

1. If $A \in N \cup \Sigma$, then A is a parse tree. Its root is A and the yield is A .
2. If $X \rightarrow \epsilon \in R$, then $X - \epsilon$ is a parse tree. Its root is X , yield is ϵ .
3. If $X \rightarrow X_1 X_2 \cdots X_n \in R$, and R_1, \dots, R_n are parse trees with roots X_1, \dots, X_n , and yield y_1, \dots, y_n , then $X - T_1, X - T_2, \dots, X - T_n$ is a parse tree with root x and yield $y_1 \cdots y_n$.

Definition 6

The definition of the language of a grammar can be equivalently stated as: $\mathcal{L}(G) = \{x \in \Sigma^* \mid x \text{ is the yield of a parse tree of } G\}$.

Note 2

Note that the yield of a parse tree is a string in $(N \cup \Sigma)^*$.

Definition 7

L is a context-free language (CFL) if there exists a grammar G such that $L = \mathcal{L}(G)$.

Example 5

Denote $S \rightarrow \epsilon \mid 0S1 \mid SS$. Then $S \Rightarrow SS \Rightarrow S0S1 \Rightarrow S01$ proves that the yield is neither the pre-order traversal or the bfs traversal.

Further, consider the string 0101.

$S \Rightarrow SS \Rightarrow 0S1S \Rightarrow 0S10S1 \Rightarrow 010S1 \Rightarrow 0101$

and

$S \Rightarrow SS \Rightarrow S0S1 \Rightarrow 0S10S1 \Rightarrow 010S1 \Rightarrow 0101$

and so on, all give the same parse tree.

So this shows that there can be multiple derivations for the same parse tree (consider labelling the expanded node by expansion time for uniqueness).

Definition 8

$Z_0 \Rightarrow X_1 \Rightarrow \cdots \Rightarrow Z_n$ is called a leftmost derivation if in each production $Z_i \Rightarrow Z_{i+1}$ a production rule is applied to the leftmost non-terminal of Z_i (essentially a pre-order traversal).

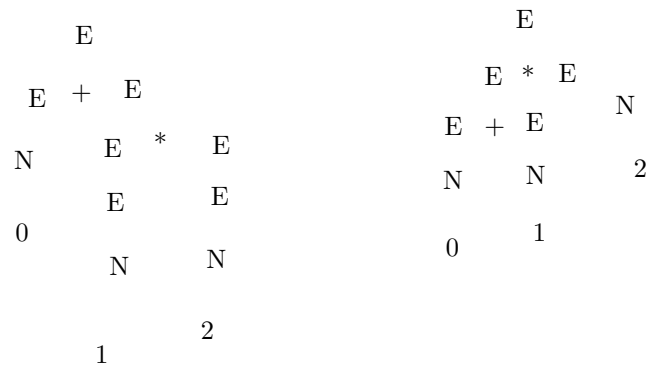
Claim 0.3

There is a one-one correspondence between parse trees yielding strings in $\mathcal{L}(G)$ and the leftmost derivations of strings in $\mathcal{L}(G)$.

Proof. Exercise – should be intuitively clear. □

We show that there can be multiple parse trees for the same string:

Consider $E \rightarrow E + E \mid E * E \mid (E) \mid N$.



Definition 9

Grammar G is said to be ambiguous if there is some string $x \in \mathcal{L}(G)$ which has two leftmost derivations (i.e., two parse trees).