

University of Waterloo

CS 241 Winter 2018

Review

Formal definition of DFA:

A DFA is a 5-tuple $(\Sigma, Q, q_0, A, \delta)$ where:

- finite alphabet Σ
- finite set of states Q
- start state $q_0 \in Q$
- set of final/accepting states $A \subseteq Q$
- transition function $\delta : Q \times \Sigma \rightarrow Q$

Formal definition of NFA:

A NFA is a 5-tuple (Σ, Q, q_0, A, T) where:

- finite alphabet Σ
- finite set of states Q
- start state $q_0 \in Q$
- set of final/accepting states $A \subseteq Q$
- transition relation $T : Q \times \Sigma \rightarrow 2^Q$

For example, if $Q = \{A, B\}$, then $2^Q = \{\{\epsilon\}, \{A\}, \{B\}, \{A, B\}\}$. $|2^Q| = 2^Q$.
NFA to DFA: subset construction

ϵ -NFA to NFA:

1. take ϵ shortcuts; replace ϵx with $x \forall x \in \Sigma$
2. pull back final states
3. remove ϵ transitions
4. remove dead states (states that have no transitions into it)

Scanning: the Simplified Maximal Munch Algorithm

1. Start at the start state of the scanning DFA.
2. Read characters until an error is reached or input is exhausted, keep tracking of the current state and the previous state. In addition, keep track of the characters read.
3. If the input is exhausted and the current state is an accepting state, emit the characters read as a token. Otherwise, signal an error in tokenizing.
4. If an error is reached and the previous state is an accepting state, emit the characters read (excluding the one that caused the error) as a token. Then resume from step 1 using the remaining input (starting at the character that caused the error) with both the current and previous states set to the start state of the scanning DFA.

Regular expressions:

RE is

- \emptyset , or
- ϵ , or
- a , where $a \in \Sigma$, or
- E_1E_2 where E_1, E_2 are REs (series), or
- $E_1|E_2$ where E_1, E_2 are REs (parallel), or
- E^* where E is RE (feedback)

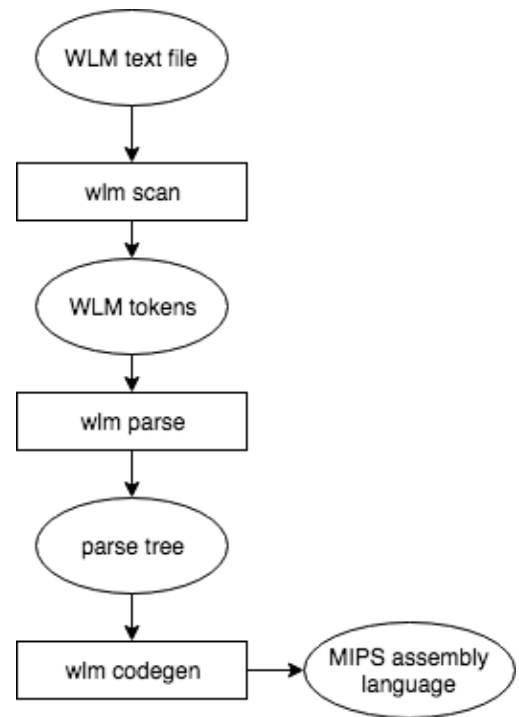
Definition of Regular Language:

A regular language is a language which is

- specified by a regular expression
- recognized by an ϵ -NFA
- recognized by an NFA
- recognized by a DFA

Compilation:

- Lexical analysis: scanning/tokenizing
- Syntactic analysis: parsing
- Context-sensitive (semantic) analysis
- Synthesis (code generation)



Context-free grammar(CFG):

G : CFG

$L(G)$: set of words specified by G (i.e. the language specified by G)

N : finite set of non-terminals (non-ending)

T : finite set of terminals (ending)

P : finite set of production rules (rewriting rules)

S : $S \in N$, start symbol

We say that $\alpha A \beta$ directly derives $\alpha \gamma \beta$ if there exists a production rule $A \rightarrow \gamma$. Also called a *derivation step*.

We say that $\alpha A \beta$ derives $\alpha \gamma \beta$ if there exist two or more derivation steps such that $\alpha A \beta \rightarrow^* \alpha \gamma \beta$.

Leftmost derivation: when there are 2 or more non-terminals in a derivation, we pick the leftmost non-terminal.

Rightmost derivation: when there are 2 or more non-terminals in a derivation, we pick the rightmost non-terminal.

Ambiguity in CFGs:

A string x is ambiguous if $x \in L(G)$ and there are more than one parse trees for x .

A CFG G is ambiguous if some word $w \in L(G)$ is ambiguous.

A grammar is ambiguous if there is a word x such that x has

1. ≥ 2 different parse trees, or
2. ≥ 2 different leftmost derivations, or
3. ≥ 2 different rightmost derivations

left recursion: leftmost symbol of the RHS is the LHS ($E \rightarrow E + B$)

right recursion: rightmost symbol of the RHS is the LHS ($E \rightarrow B + E$)

Related to associativity.

Top-down parsing with a stack:

Invariant: derivation = input already read + stack (stack is read from the top-down)

Use a rule: pop the stack which has the LHS non-terminal, push RHS in reverse onto the stack.

Accept the input when simultaneously stack and input are empty.

The oracle: $\text{Predict}(A, x) = A \rightarrow \alpha$

A is on top of the stack, and x is the first symbol of input to be read.

LL(1) Grammar:

\forall non-terminal $A \in N, x \in T, |\text{Predict}(A, x)| \leq 1$.

L: Left-to-right input

L: leftmost derivation

1: one token of “lookahead” (in terms of input)

Constructing a Predictor Table: form search trees (search until the first terminal symbol)

Algorithm:

$\alpha, \beta \in (N \cup T)^*, x, y \in T, A \in N$

```
1 Empty( $\alpha$ ) = true if  $\alpha \rightarrow^* \epsilon$ 
2 // can  $\alpha$  disappear?
3 First( $\alpha$ ) =  $\{x \mid \alpha \rightarrow^* x\beta\}$ 
4 // starting from  $\alpha$ , what can I generate as a first terminal symbol?
5 Follow( $A$ ) =  $\{y \mid S' \rightarrow^* \alpha A y \beta\}$ 
6 // starting from the start symbol, does the terminal  $y$  ever appear following the non-
   terminal  $A$ ?
7 Predict( $A, x$ ) =  $\{A \rightarrow \alpha \mid x \in \text{First}(\alpha)\} \cup \{A \rightarrow \beta \mid x \in \text{Follow}(A) \text{ and } \text{Empty}(\beta)\}$ 
```

Algorithm for parsing:

```
1 Input:  $w$ 
2 Push  $S'$ 
3 for each  $x \in w$ 
4   while(top of stack is some  $A \in N$ ){
5     pop  $A$ 
6     if Predict( $A, x$ ) =  $\{A \rightarrow \alpha\}$ 
7       push  $\alpha$ 
8     else
9       reject
10  }
11  pop  $c$ 
12  if  $c \neq x$  reject
13 end for
14 accept  $w$ 
```

Bottom-up parsing with a stack:

stack + input to be read = current derivation (stack is read from bottom to top)

shift: push

reduce($A \rightarrow ab$): pop RHS, push LHS

oracle: in the form of a DFA

LR(0) ϵ -NFA formal definition:

Given a CFG $G = (N, T, P, S)$, construct an ϵ -NFA($Q, N \cup T, q_0, F, D$) as follows:

- $Q = \{A \rightarrow \alpha \bullet \beta \mid A \rightarrow \alpha \beta \in P\}$
- $q_0 = \{S' \rightarrow \vdash \bullet S \dashv\}$

- $D[A \rightarrow \alpha \bullet X\beta, X] = \{A \rightarrow \alpha X \bullet \beta\}$
- $D[A \rightarrow \alpha \bullet X\beta, \epsilon] = \{B \rightarrow \bullet \gamma \mid B \rightarrow \gamma \in P\}$
- $F = \{S' \rightarrow \vdash S \bullet \dashv\}$

2 actions:

If you have stack K and input a , and Ka is recognized, you can **shift**.

If you have stack K and input a , and the top of K is a state containing $A \rightarrow \alpha \bullet$ and a can follow A , **reduce** $A \rightarrow \alpha \bullet$.

If more than one of these are defined, you have a **conflict**.

Building an LR(0) automaton:

An item is a production with a dot (\bullet) somewhere on the RHS (which indicates a partially completed rule).

How to construct:

- make the start state the first rule, with the dot (\bullet) in front of the leftmost symbol of the RHS
- for each state, label an arc with the symbol that follows \bullet and advance the \bullet one position to the right in the next state
- if the \bullet precedes a non-terminal, add all productions with that non-terminal on the LHS to the current state, with the \bullet in the leftmost position

Using the automaton:

- If there is a transition out of our current state on the current input, then *shift* (push) that input onto the stack
- We know we can *reduce* if the current state has only one item and the \bullet is the rightmost symbol
- To *reduce*, pop the RHS off the stack, reread the stack (from the bottom-up), follow the transition for the LHS and push the LHS onto the stack

Conflict: shift-reduce, reduce-reduce

If such conflicts are present, the grammar is not LR(0)

Adding a lookahead token to the automaton fixes the conflict.

For each $A \rightarrow \alpha$, attach $\text{Follow}(A)$.

Increasing efficiency: store (state, input) on stack, and start the transducer at the top of the stack.

Outputting a derivation (parse tree):

- Create a “tree stack”
- Each time we reduce, pop the RHS nodes from tree stack
- Push the LHS node and make its children the nodes we just popped

Traversing the tree twice:

- Symbol table: variable values
- Detect semantic errors

Context-sensitive analysis:

Input: variable names, procedure names, parse tree (syntactically valid)

Output: ERROR if there are any context-sensitive errors; (decorated) parse tree otherwise

Selected problems from tutorials:

Write context-free grammars for the following languages:

- (a) The language $L = \{\text{my, name, is, inigo, montoya}\}$.
- $S \rightarrow \text{my}$
 - $S \rightarrow \text{name}$
 - $S \rightarrow \text{is}$
 - $S \rightarrow \text{inigo}$
 - $S \rightarrow \text{montoya}$
- (b) The language of all words over $\Sigma = \{a, b, c\}$ not beginning with b .
- $S \rightarrow aA$
 - $S \rightarrow cA$
 - $A \rightarrow aA$
 - $A \rightarrow bA$
 - $A \rightarrow cA$
 - $A \rightarrow \epsilon$
 - $S \rightarrow \epsilon$
- (c) The language defined by the regular expression $(0|1)^*((00)^*1)^*010$.
- $S \rightarrow AB010$
 - $A \rightarrow 0A$
 - $A \rightarrow 1A$
 - $A \rightarrow \epsilon$
 - $B \rightarrow C1B$
 - $B \rightarrow \epsilon$
 - $C \rightarrow 00C$
 - $C \rightarrow \epsilon$

- (d) The language $L = \{a^n b^n c^m d^m : n, m \in \mathbb{N}\}$, where a^n means “ n copies of a in a row”.

$$S \rightarrow AB$$

$$A \rightarrow aAb$$

$$A \rightarrow \epsilon$$

$$B \rightarrow cBd$$

$$B \rightarrow \epsilon$$

- (e) The language of palindromes over $\Sigma = \{a, b\}$.

$$S \rightarrow a$$

$$S \rightarrow b$$

$$S \rightarrow \epsilon$$

$$S \rightarrow aSa$$

$$S \rightarrow bSb$$

$$S \rightarrow \epsilon$$

Not on final but interesting stuff:

LR(1) but not LALR(1) and not SLR(1):

$S \rightarrow AaBa$

$S \rightarrow BbAb$

$A \rightarrow c$

$B \rightarrow c$

After building the LR(0) state set for this grammar, get:

$A \rightarrow c$

$B \rightarrow c$

The SLR(1) construction uses the follow set for A , $\{a, b\}$, to determine whether to reduce $A \rightarrow c$ and the follow set for B , $\{a, b\}$, to determine whether to reduce $B \rightarrow c$. Since these sets intersect there's a reduce/reduce conflict.

LR(1) is clever enough to build two different states, one containing

$A \rightarrow c \{a\}$

$B \rightarrow c \{b\}$

and the other containing

$A \rightarrow c \{b\}$

$B \rightarrow c \{a\}$

LR(1) and LALR(1) but not SLR(1):

$S \rightarrow AaBa$

$S \rightarrow BbAb$

$S \rightarrow c$

$A \rightarrow c$

$B \rightarrow c$

For LALR(1), get states:

$A \rightarrow c$

$B \rightarrow c$

and

$A \rightarrow c$

$B \rightarrow c$

$S \rightarrow c$

If you look at the grammar carefully, you'll see that the only reductions that make sense are:

$A \rightarrow c \{b\}$

$B \rightarrow c \{a\}$

and

$A \rightarrow c \{a\}$

$B \rightarrow c \{b\}$

$S \rightarrow c \{\text{EOF}\}$

But SLR(1) is not powerful enough to make this inference as it uses the same lookahead set for each reduce action, irrespective of the state that it occurs in.