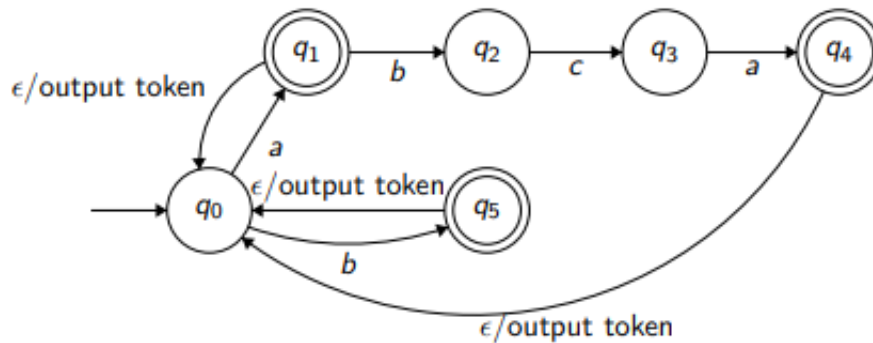


CS 241, Lecture 9 - Maximal Munch and Context Free Grammars

Tues, Feb 05, 2019

1 Maximal Munch and Simplified Maximal Munch

- The general idea is to consume the largest possible token that makes sense, then produce the token and then proceed.
- Maximal munch differs from simplified maximal munch:
 - Maximal munch consumes chars until there is no valid transition. If there are still characters, backtrack to the last valid accepting state and resume.
 - Simplified maximal munch consumes chars until there is no longer a valid transition. If the state is accepting, produce the token and proceed. Otherwise go to error state.
- Observe the following DFA for $w = \text{“ababca”}$:
 $\Sigma = \{a, b, c\}$, $L = \{a, b, abca\}$, consider $w = ababca$.
Note that $w \in LL^*$. What follows is an ϵ -NFA for LL^* based on our algorithm:



- With SMM, when we hit the second “a”, we stop and output an error, which is not the correct answer.
- With MM, the second “a” is reached, and the last accepting state was “a”. We backtrack back to “a” and resume munch (go back to q_0 , resume by consuming “b”. We would have to keep track of the last accepting state.
- Note SMM is usually good enough and thus we typically use this without issue - that is, we make our language work with SMM rather than the other way around.

- Examples of the MM and SMM algorithms:

Algorithm for Maximal Munch

Algorithm 1 Maximal Munch

```

1:  $s = q_0, t_a = \epsilon, t_{cur} = \epsilon$ 
2:  $str = input, pos = 0, posAccepting = -1$ 
3: while  $pos \neq len(str)$  do
4:    $c = str[pos], s = \delta(s, c), t_{cur} = t_{cur} + c$ 
5:   if  $s == ERROR$  then
6:     if  $t_a = \epsilon$  then
7:       Fatal Error
8:     end if
9:     Output  $t_a$ 
10:     $s = q_0, t_a = \epsilon, t_{cur} = \epsilon, pos = posAccepting$ 
11:  else if  $s \in A$  then
12:     $posAccepting = pos$ 
13:     $t_a = t_{cur}$ 
14:  end if
15:   $pos = pos + 1$ 
16: end while
17: if  $s \in A$  then
18:   Output  $t$  and Accept
19: else
20:   Reject/Crash/Fatal Error
21: end if

```

Algorithm for Simplified Maximal Munch

Algorithm 2 Simplified Maximal Munch

```

1:  $s = q_0$ 
2:  $t_a = \epsilon$ 
3: while not EOF do
4:    $c = read\_char()$ 
5:   if  $\delta(s, c) == ERROR$  then
6:     if  $s \in A$  then
7:       Output  $t_a$ 
8:        $s = q_0, t_a = \epsilon$ 
9:     else
10:      Reject/Crash/Fatal Error
11:    end if
12:  else
13:     $s = \delta(s, c)$ 
14:     $t_a = t_a + c$ 
15:  end if
16: end while
17: if  $s \in A$  then
18:   Output  $t_a$  and Accept
19: else
20:   Reject/Crash/Fatal Error
21: end if

```

- This concludes the process of scanning.

2 Syntactic Analysis - Context-Free Grammars

- Things we have to consider: syntax (is the order of tokens correct, are parentheses balanced?) and semantics (does what is written make sense?).
- A **grammar** is the language of languages - they help us describe what we are and are not allowed to say.
- **Context-free grammars** is a 4-tuple (N, Σ, P, S) where:
 - N is a finite non-empty set of non-terminal symbols (symbols you cannot stop on)
 - Σ is an alphabet, or in other words a set of non-empty terminal symbols, and $N \cap \Sigma = \emptyset$.
 - P is a finite set of productions, each of the form $A \rightarrow \beta$ where $A \in N$ and $\beta \in (N \cup \Sigma)^*$.
 - $S \in N$ is a starting symbol.
- We set $V = N \cup \Sigma$ to denote the vocabulary - the set of *all* symbols in our language.
- For example, in rustcc, we defined various CFGs, such as Fn containing BlockItems which contained Statements or Declarations which contained...
- **Conventions:**
 - Lower case letters from the start of the alphabet are elements of Σ
 - Lower case letters from the end of the alphabet are elements of Σ^* (words)
 - Upper case letters from the start of the alphabet are elements of N
 - S is always our start symbol
 - Greek letters like α, β, γ are elements of $V^* = (N \cup \Sigma)^*$.
- For example, consider $\Sigma = \{ (,) \}$, and let $L = \{ w : w \text{ is a balance string of parentheses} \}$. Thus, $S \rightarrow \epsilon, S \rightarrow (S), S \rightarrow SS \Rightarrow S \rightarrow \epsilon|(S)|SS$
- A **derivation** over a CFG (N, Σ, P, S) is such that:
 - α derives β and we write $\alpha \Rightarrow \beta$ iff β can be obtained from α using a rule from P .
 - $\alpha A \beta \Rightarrow \alpha \gamma \beta$ iff there is a rule $A \rightarrow \gamma$ in P .
 - $\alpha \Rightarrow^* \beta$ iff a derivation exists, that is, there exists $\delta_i \in V^*$ for $0 \leq i \leq k$ such that $\alpha = \delta_0 \Rightarrow \delta_1 \Rightarrow \dots \Rightarrow \delta_k = \beta$. Note k can be 0.
- Another example: find a derivation of $((()))$. Recall our above CFG:

$$\begin{aligned}
 S &\Rightarrow (S) \Rightarrow (SS) \Rightarrow ((S)S) \\
 &\Rightarrow ((S)(S)) \Rightarrow ((\epsilon)(S)) \\
 &\Rightarrow (()()) \\
 S &\Rightarrow^* (()()) \quad , \text{ short form for the above implications}
 \end{aligned}$$

- Why is it “context-free”? It’s because our grammar does not care for the context - that is, it does not care for *where* your symbols are.
- This is the opposite of **context-bounded grammars**, which is where the context of the other symbols around other symbols *will* affect our productions (note we don’t *really* need to know about this)
- We define the **language** of a $\text{CFG}(N, \Sigma, P, S)$ to be $L(G) = \{w \in \Sigma^* : S \Rightarrow^* w\}$.
- A language is **context-free** iff there exists a CFG G such that $L = L(G)$.
- Informally, we can show regular languages are context-free:
 1. \emptyset : $(\{S\}, \{a\}, \emptyset, S)$
 2. $\{\epsilon\}$: $(\{S\}, \{a\}, S \rightarrow \epsilon, S)$.
 3. $\{a\}$: $(\{S\}, \{a\}, S \rightarrow a, S)$.
 4. Union: $\{a\} \cup \{b\}$: $(\{S\}, \{a, b\}, S \rightarrow a|b, S)$.
 5. Concatenation: $\{ab\}$: $(\{S\}, \{a, b\}, S \rightarrow ab, S)$.
 6. Kleene Star: $\{a\}^*$: $(\{S\}, \{a\}, S \rightarrow Sa|\epsilon, S)$.