

CSC301 Lecture 10

1910456 Mir Shafayat Ahmed

August 26, 2021

Regex to CFG

We use the following replacements when converting different operations in regex to CFG.

$\emptyset = \text{grammar with no rules}$

$\varepsilon = G \rightarrow \varepsilon$

$A^* = G \rightarrow GA \mid \varepsilon$

$A + B = G \rightarrow A \mid B$

$AB = G \rightarrow AB$

So, to convert $(0+1)^*111$, We can divide it into two parts, one with the $()^*$ and one with the 111 and use the other replacements to find,

$S \rightarrow N111$

$N \rightarrow N0 \mid N1 \mid \varepsilon$

Ambiguity

If a string has more than one parse tree, we call the CFG ambiguous. Such as for the string $1+2*2$, without the parenthesis, the CFG would be ambiguous.

For a CFG that has a rule $N \rightarrow NN$ compared to a rule that says $N \rightarrow DN$, there is ambiguity on where to add the next sequence. (To add to the left or to the right).

$S \rightarrow SS \mid x$

But we can fix CFG of Figure 1 by simply changing it to

$S \rightarrow Sx \mid x$

And now, subsequent additions to the tree can only take one path.



Figure 1: Different parse tree for same string

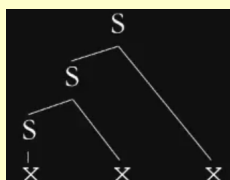


Figure 2: New Unique Parse Tree

Disambiguation

When we are looking to make a language unambiguous, we might not succeed as some languages are inherently unambiguous. Also, there is no general procedure to follow.

In programming languages, ambiguity comes from not using parenthesis. In spoken languages as well, there can be some ambiguity.

"She fixed the light on the top floor"

Parsing

For the CFG,

$$S \rightarrow 0S1 \mid 1S0S \mid T$$

$$T \rightarrow S \mid \varepsilon$$

If we want to build a parse tree for the string '0011', we can go on and on without any final destination.

It is because of a few reasons,

- The string might not be in the language and so we are on a never-ending search.
- For longer strings, the permutation of parse trees would be too many.

We therefore have an idea to:

Stop when the length of a derived string exceeds the required string

But here we encounter two problems:

1. Due to empty transitions, strings containing Variables may turn into empty strings and reduce the length of the string later on.
2. Rules that have Unit transitions like $A \rightarrow B$ may continue looping forever.

So, to resolve the two issues, we modify our CFG so that there are no ε and unit productions. After which, **we can be sure that any derivation can only grow the string and not shrink or remain constant in length.**

Removing ε -productions

If the start variable is nullable,

$$S \rightarrow \varepsilon$$

we add a new start variable S' and add a new production,

$$S' \rightarrow S \mid \varepsilon$$

For other nullable variables,

1. We identify a nullable variable and remove the ε -production
2. We look for references of that variable in other productions and replace that reference with what would've happened if the nullable variable wasn't removed.

Example

For the CFG,

$$\begin{aligned}
 S &\rightarrow ACD \\
 A &\rightarrow a \\
 B &\rightarrow \varepsilon \\
 C &\rightarrow ED \mid \varepsilon \\
 D &\rightarrow BC \mid b \\
 E &\rightarrow b
 \end{aligned} \tag{1}$$

We identify B as a nullable variable, so we remove $B \rightarrow \varepsilon$. Now we look for references and see that B is referenced in $D \rightarrow BC$. So we replace B with what it was supposed to be, i.e 'empty'. And so, we get,

$$\begin{aligned}
 S &\rightarrow ACD \\
 A &\rightarrow a \\
 C &\rightarrow ED \mid \varepsilon \\
 D &\rightarrow C \mid b \\
 E &\rightarrow b
 \end{aligned} \tag{2}$$

Now we look for more nullable variables and see that C is nullable. So we remove $C \rightarrow \varepsilon$. We see that C is referenced in the beginning, so we change it to $S \rightarrow AD$. But, now we cannot form something like 'AEDD' that was possible before, therefore we keep $S \rightarrow ACD$ as well. D also references C, so we replace it with $D \rightarrow \varepsilon$. But again, D could be ED before, so we keep $D \rightarrow C$

$$\begin{aligned}
S &\rightarrow ACD \\
S &\rightarrow AD \\
A &\rightarrow a \\
C &\rightarrow ED \\
D &\rightarrow C \mid b \mid \varepsilon \\
E &\rightarrow b
\end{aligned} \tag{3}$$

We repeat and remove $D \rightarrow \varepsilon$. Now $S \rightarrow ACD$ can be $S \rightarrow AC$ as well. For the other reference, the same logic applies and we end up with,

$$\begin{aligned}
S &\rightarrow ACD \mid AC \\
S &\rightarrow AD \mid A \\
A &\rightarrow a \\
C &\rightarrow ED \mid E \\
D &\rightarrow C \mid b \\
E &\rightarrow b
\end{aligned} \tag{4}$$

Finally we see that there aren't any nullable variables anymore. So we have solved our first problem regarding disambiguation.