CS164 - Left-recursion Elimination example (Weimer)

Consider the grammar:

$$
\begin{aligned}
A &\rightarrow B \mid a \mid CBD \\
B &\rightarrow C \mid b \\
C &\rightarrow A \mid c \\
D &\rightarrow d
\end{aligned}
$$

Some strings in the language of this grammar are: $a$, $cbd$, etc. Notice that this grammar is not **immediately** left-recursive in that there is no single production $X \rightarrow X\alpha$. However, it is left-recursive because there are valid derivations of the form $A \rightarrow^* A\alpha$ (and $B \rightarrow^* B\beta$ and $C \rightarrow^* C\delta$). Let's demonstrate one: $A \rightarrow B \rightarrow C \rightarrow A$, so $A \rightarrow^* A$.

To warm up, let's compute $First()$ and $Follow()$ sets for this grammar.

$First(A)$ must contain $First(B)$, $\{a\}$ and $First(CBD)$. $First(B)$ must contain $First(C)$ and $\{b\}$. $First(C)$ must contain $First(A)$ and $\{c\}$. $First(CBD) = First(C)$ as before. So $First(A) = \{a, b, c\}$. Similarly, $First(B) = First(C) \bigcup \{b\} = \{a, b, c\}$ and $First(C) = First(A) \bigcup \{c\} = \{a, b, c\}$. $First(D) = \{d\}$.

To compute $Follow(A)$ we look for every occurrence of $A$ on the right-hand side of a production. We find one in $C \rightarrow A$, but it that $A$ is at the direct end of the production, we get that $Follow(A)$ includes $Follow(C)$. Now we look for $C$ on the right-hand side of a production and find it in $A \rightarrow CBD$. This time there is something to the right of $C$, so we get that $Follow(C)$ contains $First(B)$. However, we also see a $C$ on the right of $B \rightarrow C$, so $Follow(C)$ contains $Follow(B)$. Looking for $B$'s on the right we find $A \rightarrow CBD$, so $Follow(B)$ contains $First(D) = \{d\}$l. So $Follow(A) = Follow(B) = Follow(C) = \{a, b, c, d\}$. Since $D$ appears in the production $A \rightarrow CBD$, we have that $Follow(D)$ includes $Follow(A)$, so $Follow(D) = \{a, b, c, d\}$ as well.

OK, messy grammar. Now let's eliminate left-recursion. The first step is to make all left-recursion **immediate** by doing some substitutions. For example, since we have $A \rightarrow B \rightarrow C \rightarrow A$, we need to take the production $A \rightarrow B$ and replace it with $A \rightarrow C$ and $A \rightarrow b$. That gives us:

$$
A \rightarrow C \mid b \mid a \mid CBD
$$

But we're not done, since we can have $A \rightarrow C \rightarrow A$. So now we need to substitute in for $C$ in that production. Let's do that once:

$$
A \rightarrow A \mid c \mid b \mid a \mid CBD
$$

To get here, we just removed the production $A \rightarrow C$ and replaced it by $A \rightarrow A$ and $A \rightarrow c$ (we got those two right-hand sides from the productions $C \rightarrow A$ and $C \rightarrow c$). Now we're almost done, just one more possible non-immediate left-recursion. Let's substitute it away:

$$
A \rightarrow A \mid c \mid b \mid a \mid ABD \mid cBD
$$

Huzzah! Now $A$ has only immediate left-recursion. And actually, there is no other left-recursion left in the grammar now, since we can no longer derive $B \rightarrow^* B\beta$ or $C \rightarrow^* C\delta$. So we can leave the $B \rightarrow$ and $C \rightarrow$ and $D \rightarrow$ productions alone and concentrate on eliminating left-recursion from $A$.

First, let's group the productions into those that are left-recursive and those that are not:

$$
\begin{aligned}
A &\rightarrow A \mid ABD \\
&\mid c \mid b \mid a \mid cBD
\end{aligned}
$$

Now imagine that you actually have this grammar before you. You can expand things for a long time by just using the first to productions: $A \rightarrow ABD \rightarrow ABDBD \rightarrow ABDBD$, etc. But eventually you have to settle down and use

one of the other productions: $A \to ABD \to ABDBD \to cBDBD$ and then the chain stops. So we get the idea that $A$ can eventually produce something that starts with $c$, $b$, $a$ or $cBD$ and ends with a list of $BD$'s.

One other thing to note is that the production $A \to A$ itself is useless – it does not change the language of the grammar and can be safely dropped. So here's the revised left-recursive grammar:

$$
\begin{aligned}
A \quad &\to \quad ABD \\
&| \quad c \,|\, b \,|\, a \,|\, cBD
\end{aligned}
$$

Now let's break that down into $A$ and $A'$. We reasoned above that an $A$ goes to $c|b|a|cBD$ followed by a list of $BD$'s. Let's make the first bit the $A$ and make the list of $BD$'s the $A'$.

$$
\begin{aligned}
A \quad &\to \quad cA' \,|\, bA' \,|\, aA' \,|\, cBDA' \\
A' \quad &\to \quad \varepsilon \,|\, BDA'
\end{aligned}
$$

We're done. You can check and see that every production for $A$ is of the form $A \to \alpha A'$ and that $A'$ really does define a (possibly empty) list of $BD$'s. Let's do one more example, just for fun. Consider the grammar:

$$
\begin{aligned}
Q \quad &\to \quad QED \,|\, q \\
E \quad &\to \quad e \\
D \quad &\to \quad NFA \,|\, d \\
N \quad &\to \quad DFA \,|\, n \\
F \quad &\to \quad f \\
A \quad &\to \quad a
\end{aligned}
$$

This grammar is left recursive. In fact, it is immediately left-recursive in one place and non-immediately left-recursive in two places. First, let's substitute to get rid of the non-immediate left recursion. Consider the derivation: $D \to NFA \to DFAFA$. Since it ends up being left recursive, we must substitute. Take the production $D \to NFA$ and remove it. Then for every production $N \to \alpha_i$, add a production $D \to \alpha_i FA$. That gives us:

$$
D \quad \to \quad DFAFA \,|\, nFA \,|\, d
$$

Notice again that in one fell swoop we have eliminated the whole chain of non-immediate left-recursion: we can no longer derive $N \to^* N\beta$. So now our grammar looks like:

$$
\begin{aligned}
Q \quad &\to \quad QED \,|\, q \\
E \quad &\to \quad e \\
D \quad &\to \quad DFAFA \,|\, nFA \,|\, d \\
N \quad &\to \quad DFA \,|\, n \\
F \quad &\to \quad f \\
A \quad &\to \quad a
\end{aligned}
$$

Now it's time to eliminate the immediate left recursion. Let's start with $Q \to QED|q$. Once again, we can derive strings like $Q \to QEDEDED$, but eventually we have to stop and use $Q \to q$. Taking $Q$ to be the $q$ bit and $Q'$ to be the list of $ED$'s, we get:

$$
\begin{aligned}
Q \quad &\to \quad qQ' \\
Q' \quad &\to \quad \varepsilon \,|\, EDQ'
\end{aligned}
$$

2

Now let's look at $D \rightarrow DFAFA|nFA|d$. We see that we can make a huge list of $FAFA$' s using the first production but we eventually have to start with $nFA$ or $d$. Let's make $D'$ the list and $D$ the first bit.

$$
\begin{aligned}
D &\rightarrow nFAD' \,|\, dD' \\
D' &\rightarrow \varepsilon \,|\, FAFAD'
\end{aligned}
$$

OK, that's all the left-recursion. The final grammar is:

$$
\begin{aligned}
Q &\rightarrow qQ' \\
Q' &\rightarrow \varepsilon \,|\, EDQ' \\
E &\rightarrow e \\
D &\rightarrow nFAD' \,|\, dD' \\
D' &\rightarrow \varepsilon \,|\, FAFAD' \\
N &\rightarrow DFA \,|\, n \\
F &\rightarrow f \\
A &\rightarrow a
\end{aligned}
$$

Huzzah, we are done.