



CSE 322

Construction

of Reduced Grammars and Elimination of null and unit productions

Lecture #27

How to “simplify” CFGs?

Three ways to simplify/clean a CFG

(clean)

1. Eliminate *useless symbols*

(simplify)

2. Eliminate ϵ -productions

$A \Rightarrow \epsilon$

3. Eliminate *unit productions*

$A \Rightarrow B$

Eliminating useless symbols

Eliminating *useless symbols*

A symbol X is reachable if there exists:

- $S \rightarrow^* \alpha X \beta$

A symbol X is generating if there exists:

- $X \rightarrow^* w$,
 - for some $w \in T^*$

For a symbol X to be “useful”, it has to be both reachable *and* generating

- $S \rightarrow^* \alpha X \beta \rightarrow^* w'$, for some $w' \in T^*$

reachable generating

Algorithm to detect useless symbols

1. First, eliminate all symbols that are *not* generating
2. Next, eliminate all symbols that are *not* reachable

Is the order of these steps important,
or can we switch?

Example: Useless symbols

- $S \rightarrow AB \mid a$
 - $A \rightarrow b$
1. A, S are generating
 2. B is *not generating* (and therefore B is useless)
 3. \Rightarrow Eliminating $B...$ (i.e., remove all productions that involve B)
 1. $S \rightarrow a$
 2. $A \rightarrow b$
 4. Now, A is *not reachable* and therefore is useless
 5. Simplified G :
 1. $S \rightarrow a$

What would happen if you reverse the order:
i.e., test reachability before generating?

Will fail to remove: $A \rightarrow b$

$$X \rightarrow^* w$$

Algorithm to find all generating symbols

- Given: $G=(V,T,P,S)$
- Basis:
 - Every symbol in T is obviously generating.
- Induction:
 - Suppose for a production $A \rightarrow \alpha$, where α is generating
 - Then, A is also generating

Algorithm to find all reachable symbols

$$S \rightarrow^* \alpha X \beta$$

- Given: $G=(V,T,P,S)$
- Basis:
 - S is obviously reachable (from itself)
- Induction:
 - Suppose for a production $A \rightarrow \alpha_1 \alpha_2 \dots \alpha_k$, where A is reachable
 - Then, all symbols on the right hand side, $\{\alpha_1, \alpha_2, \dots, \alpha_k\}$ are also reachable.

Eliminating ϵ -productions

$A \Rightarrow \epsilon$

X

Eliminating ε -productions

$$A \rightarrow \varepsilon$$

Caveat: It is *not* possible to eliminate ε -productions for languages which include ε in their word set

So we will target the grammar for the rest of the language

Theorem: If $G=(V,T,P,S)$ is a CFG for a language L , then $L \setminus \{\varepsilon\}$ has a CFG without ε -productions

Definition: A is “nullable” if $A \rightarrow^* \varepsilon$

- If A is nullable, then any production of the form “ $B \rightarrow CAD$ ” can be simulated by:
 - $B \rightarrow CD \mid CAD$
 - This can allow us to remove ε transitions for A

Algorithm to detect all nullable variables

- Basis:
 - If $A \rightarrow \varepsilon$ is a production in G , then A is nullable (note: A can still have other productions)
- Induction:
 - If there is a production $B \rightarrow C_1C_2...C_k$, where *every* C_i is nullable, then B is also nullable

Eliminating ϵ -productions

Given: $G=(V,T,P,S)$

Algorithm:

1. Detect all nullable variables in G
2. Then construct $G_1=(V,T,P_1,S)$ as follows:
 - i. For each production of the form: $A \rightarrow X_1 X_2 \dots X_k$, where $k \geq 1$, suppose m out of the k X_i 's are nullable symbols
 - ii. Then G_1 will have 2^m versions for this production
 - i. i.e, all combinations where each X_i is either present or absent
 - iii. Alternatively, if a production is of the form: $A \rightarrow \epsilon$, then remove it

Example: Eliminating ϵ -productions

- Let L be the language represented by the following CFG G :

- $S \rightarrow AB$
- $A \rightarrow aAA \mid \epsilon$
- $B \rightarrow bBB \mid \epsilon$

Goal: To construct G_1 , which is the grammar for $L - \{\epsilon\}$

- Nullable symbols: $\{A, B\}$
- G_1 can be constructed from G as follows:
 - $B \rightarrow b \mid bB \mid bB \mid bBB$
 - $\Rightarrow B \rightarrow b \mid bB \mid bBB$
 - Similarly, $A \rightarrow a \mid aA \mid aAA$
 - Similarly, $S \rightarrow A \mid B \mid AB$

- Note: $L(G) = L(G_1) \cup \{\epsilon\}$

Simplified
grammar

G_1 :

- $S \rightarrow A \mid B \mid AB$
- $A \rightarrow a \mid aA \mid aAA$
- $B \rightarrow b \mid bB \mid bBB$

+

- $S \rightarrow \epsilon$

Eliminating unit productions

$A \Rightarrow B$



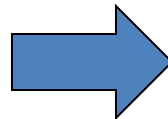
← B has to be a variable

What's the point of removing unit transitions ?

Will save #substitutions

E.g.,

$A \Rightarrow B \mid \dots$
 $B \Rightarrow C \mid \dots$
 $C \Rightarrow D \mid \dots$
 $D \Rightarrow xxx \mid yyy \mid zzz$



$A \Rightarrow xxx \mid yyy \mid zzz \mid \dots$
 $B \Rightarrow xxx \mid yyy \mid zzz \mid \dots$
 $C \Rightarrow xxx \mid yyy \mid zzz \mid \dots$
 $D \Rightarrow xxx \mid yyy \mid zzz$

before

after

Eliminating unit productions

$$A \rightarrow B$$

- Unit production is one which is of the form $A \rightarrow B$, where both A & B are variables
- E.g.,
 1. $E \rightarrow T \mid E+T$
 2. $T \rightarrow F \mid T * F$
 3. $F \rightarrow I \mid (E)$
 4. $I \rightarrow a \mid b \mid Ia \mid Ib \mid I0 \mid I1$
- How to eliminate unit productions?
 - Replace $E \rightarrow T$ with $E \rightarrow F \mid T * F$
 - Then, upon recursive application wherever there is a unit production:
 - $E \rightarrow F \mid T * F \mid E+T$ (substituting for T)
 - $E \rightarrow I \mid (E) \mid T * F \mid E+T$ (substituting for F)
 - $E \rightarrow a \mid b \mid Ia \mid Ib \mid I0 \mid I1 \mid (E) \mid T * F \mid E+T$ (substituting for I)
 - Now, E has no unit productions
 - Similarly, eliminate for the remainder of the unit productions

The Unit Pair Algorithm: to remove unit productions

- Suppose $A \rightarrow B_1 \rightarrow B_2 \rightarrow \dots \rightarrow B_n \rightarrow \alpha$
- Action: Replace all intermediate productions to produce α directly
 - i.e., $A \rightarrow \alpha; B_1 \rightarrow \alpha; \dots B_n \rightarrow \alpha;$

Definition: (A, B) to be a “*unit pair*” if $A \rightarrow^* B$

- We can find all unit pairs inductively:
 - Basis: Every pair (A, A) is a unit pair (by definition). Similarly, if $A \rightarrow B$ is a production, then (A, B) is a unit pair.
 - Induction: If (A, B) and (B, C) are unit pairs, and $A \rightarrow C$ is also a unit pair.

The Unit Pair Algorithm: to remove unit productions

Input: $G=(V,T,P,S)$

Goal: to build $G_1=(V,T,P_1,S)$ devoid of unit productions

Algorithm:

1. Find all unit pairs in G
2. For each unit pair (A,B) in G :
 1. Add to P_1 a new production $A \rightarrow \alpha$, for every $B \rightarrow \alpha$ which is a *non-unit* production
 2. If a resulting production is already there in P , then there is no need to add it.

Example: eliminating unit productions

Unit pairs		Only non-unit productions to be added to P_1
$G:$ 1. $E \rightarrow T \mid E+T$ 2. $T \rightarrow F \mid T^*F$ 3. $F \rightarrow (\mid (E)$ 4. $I \rightarrow a \mid b \mid la \mid lb \mid l0 \mid l1$	(E,E)	$E \rightarrow E+T$
	(E,T)	$E \rightarrow T^*F$
	(E,F)	$E \rightarrow (E)$
	(E,I)	$E \rightarrow a \mid b \mid la \mid lb \mid l0 \mid l1$
	(T,T)	$T \rightarrow T^*F$
	(T,F)	$T \rightarrow (E)$
	(T,I)	$T \rightarrow a \mid b \mid la \mid lb \mid l0 \mid l1$
	(F,F)	$F \rightarrow (E)$
	(F,I)	$F \rightarrow a \mid b \mid la \mid lb \mid l0 \mid l1$
	(I,I)	$I \rightarrow a \mid b \mid la \mid lb \mid l0 \mid l1$

 $G_1:$

1. $E \rightarrow E+T \mid T^*F \mid (E) \mid a \mid b \mid la \mid lb \mid l0 \mid l1$
2. $T \rightarrow T^*F \mid (E) \mid a \mid b \mid la \mid lb \mid l0 \mid l1$
3. $F \rightarrow (E) \mid a \mid b \mid la \mid lb \mid l0 \mid l1$
4. $I \rightarrow a \mid b \mid la \mid lb \mid l0 \mid l1$

Putting all this together...

- Theorem: If G is a CFG for a language that contains at least one string other than ε , then there is another CFG G_1 , such that $L(G_1) = L(G) - \varepsilon$, and G_1 has:
 - no ε -productions
 - no unit productions
 - no useless symbols
- Algorithm:
 - Step 1) eliminate ε -productions
 - Step 2) eliminate unit productions
 - Step 3) eliminate useless symbols

Again,
the order is
important!

Why?

Consider the grammar G whose productions are $S \rightarrow aS \mid AB$, $A \rightarrow \Lambda$, $B \rightarrow \Lambda$, $D \rightarrow b$. Construct a grammar G_1 without null productions generating $L(G) - \{\Lambda\}$.



Step 1 *Construction of the set W of all nullable variables:*

$$\begin{aligned} W_1 &= \{A_1 \in V_N \mid A_1 \rightarrow A \text{ is a production in } G\} \\ &= \{A, B\} \end{aligned}$$

$$\begin{aligned} W_2 &= \{A, B\} \cup \{S\} \text{ as } S \rightarrow AB \text{ is a production with } AB \in W_1^* \\ &= \{S, A, B\} \end{aligned}$$

$$W_3 = W_2 \cup \emptyset = W_2$$

Thus.

$$W = W_2 = \{S, A, B\}$$

Step 2 *Construction of P' :*

- (i) $D \rightarrow b$ is included in P' .
- (ii) $S \rightarrow aS$ gives rise to $S \rightarrow aS$ and $S \rightarrow a$.
- (iii) $S \rightarrow AB$ gives rise to $S \rightarrow AB$, $S \rightarrow A$ and $S \rightarrow B$.

(**Note:** We cannot erase both the nullable variables A and B in $S \rightarrow AB$ as we will get $S \rightarrow \Lambda$ in that case.)

Hence the required grammar without null productions is

$$G_1 = (\{S, A, B, D\}, \{a, b\}, P, S)$$

where P' consists of

$$D \rightarrow b, S \rightarrow aS, S \rightarrow AB, S \rightarrow a, S \rightarrow A, S \rightarrow B$$