

Week 8: Lecture Notes

Context free Grammars (CFG)

A context free grammar is $G = (V, T, P, S)$ where

- V is a finite set (non-empty) whose elements are called variables.
- T is a finite set (non-empty) whose elements are called terminals.
- $V \cap T = \emptyset$
- S is a special variable called the start variable.
- P is a finite set of productions or rules whose elements are $A \rightarrow \alpha$ where A is a variable and α is a string of symbols from $(V \cup T)^*$

Example:

$$G = (\{E\}, \{+, *, (,), id\}, P, E)$$

where P represents the following set of rules

$$E \rightarrow E + E \quad E \rightarrow (E)$$

$$E \rightarrow E * E \quad E \rightarrow id.$$

Applying productions repeatedly, we can obtain more and more complicated expressions.

\Rightarrow denotes act of deriving
 $(id + id)^* id$ consists of only terminal symbols, it is a word in the language generated by G .

$$\begin{aligned}
 & E \Rightarrow E + E \\
 & \Rightarrow (E) * E \\
 & \Rightarrow (E) * id \\
 & \Rightarrow (E + E) * id \\
 & \Rightarrow (E + id) * id \\
 & \Rightarrow (id + id) * id
 \end{aligned}$$

Example:

L_{pal} : languages of palindromes of 0's and 1's
- not regular

- $\epsilon, 0, 1$ are palindromes
- if w is a palindrome, so are $0w0$ and $1w1$

Such recursively defined languages (L_{pal}) can be expressed formally as context free grammar (for palindromes) G_{pal}

$$G_{\text{pal}} = (\{\mathcal{E}\}, \{0, 1\}, P, \mathcal{E})$$

where P represents the set of five productions

$$\begin{array}{ll} \mathcal{E} \rightarrow \mathcal{E} & \mathcal{E} \rightarrow 0EO \\ \mathcal{E} \rightarrow 1 & \mathcal{E} \rightarrow 1EI \\ \mathcal{E} \rightarrow 0 & \end{array}$$

- The set of productions is the kernel of grammar and language specifications
- Note (Regarding production rules)
 - (i) Reverse substitution is not permitted
e.g. if $S \rightarrow AB$ is a production, then we can replace S by AB but we cannot replace AB by S
 - (ii) No inversion operation is permitted
e.g. if $S \rightarrow AB$ is a production, then it is not necessary that $AB \rightarrow S$ is a production-

Notations

1. The capita letters A, B, C, D, E and S denote variables. S is the start symbol
2. The lowercase letters a, b, c, d, e , digits, bold face strings are terminals.
3. The capital letters X, Y, Z denote symbols that may be either terminals or variables.
4. The lowercase letters u, v, w, x, y, z denote strings of terminals.
5. The lower case Greek letters α, β, γ denote strings of variables and terminals $(VUT)^*$

A CFG is often presented by simply listing its productions.

If $A \rightarrow \alpha_1, A \rightarrow \alpha_2, \dots, A \rightarrow \alpha_k$ are productions for variable , then we may express them by the notation

$A \rightarrow \alpha_1 | \alpha_2 | \dots | \alpha_k$ ' | ' stands for 'or'

Example:

$E \rightarrow E + E, E \rightarrow E * E, E \rightarrow (E), E \rightarrow id$

can be written as

$E \rightarrow E + E | E * E | (E) | id$

Derivations and languages

- $G = (V, T, P, S) \rightarrow \text{CFG}$
- \xrightarrow{G} \rightarrow relation on $(VUT)^*$

defined as follows:

if $A \rightarrow \beta$ is a production of P and $\alpha, \gamma \in (VUT)^*$

then $\alpha A \gamma \xrightarrow{G} \alpha \beta \gamma$

$\alpha A \gamma$ directly derives $\alpha \beta \gamma$ in grammar G .

- Two strings are related by \xrightarrow{G} exactly when the second is obtained from the first by one application of some production

- \xrightarrow{G}^* relation on $(VUT)^*$

defined as follows:

Suppose that $\alpha_1, \alpha_2, \dots, \alpha_m \in (VUT)^*$, $m \geq 1$ and

$\alpha_1 \xrightarrow{G} \alpha_2, \alpha_2 \xrightarrow{G} \alpha_3, \dots, \alpha_{m-1} \xrightarrow{G} \alpha_m$

Then $\alpha_1 \xrightarrow{G}^* \alpha_m$ or α_1 derives α_m in grammar G

- \xrightarrow{G}^* is reflexive, transitive closure of \xrightarrow{G}

- $\alpha \xrightarrow{G}^* \alpha$ for each string $\alpha \in (VUT)^*$

- When it is clear which grammar G is involved, we write
 \Rightarrow instead of \xrightarrow{G} and $\xrightarrow{*}$ instead of \xrightarrow{G}^*
- $\alpha \xrightarrow{i} \beta$ if α derives β by exactly i steps
- The language generated by G is
 $L(G) = \{w \mid w \in T^* \text{ and } S \xrightarrow{G}^* w\}$
i.e. a string is in $L(G)$ if
 - The string consists of solely of terminals
 - The string can be derived from S

Context free language (CFL)

L is a context free language (CFL) if $L = L(G)$ for some GFG G

Sentential form

A string $\alpha \in (VUT)^*$ is called a sentential form if $S \xrightarrow{*} \alpha$

Equivalent grammar

G_1, G_2 are equivalent if $L(G_1) = L(G_2)$

Example:

Consider CFG $G = (V, T, P, S)$ where

$V = \{S\}$, $T = \{a, b\}$, $P = \{S \rightarrow aSb, S \rightarrow ab\}$. Find $L(G)$

$$L(G) = \{w \mid w \in T^* \text{ and } S \xrightarrow[G]{*} w\}$$

$$S \Rightarrow aSb \Rightarrow aaSbb \Rightarrow a^3Sb^3 \Rightarrow \dots \Rightarrow a^{n-1}Sb^{n-1}$$

(applying $S \rightarrow aSb$ $n-1$ times)

$$\Rightarrow a^n b^n \quad (\text{applying } S \rightarrow ab)$$

$$\text{Thus, } S \xrightarrow{*} a^n b^n \in T^* \Rightarrow a^n b^n \in L(G) \text{ for } n \geq 1$$

$$\text{Now, we need to show that } L(G) \subseteq \{a^n b^n \mid n \geq 1\}$$

Let $w \in L(G)$.

Let $w \in L(G)$. w is a terminal string i.e. $w \in T^*$

The only way w can be derived from start symbol S is

$$S \xrightarrow{*} a^{n-1} S b^{n-1} \xrightarrow{\substack{\downarrow \\ \text{applying } S \rightarrow aSb \\ \text{few no. of times}}} a^n b^n \quad \text{for some } n \geq 1$$

$\xrightarrow{\substack{\downarrow \\ \text{one application of} \\ S \rightarrow ab}}$
number of S 's is reduced by 1 in the sentential form

$$\text{i.e. } L(G) \subseteq \{a^n b^n \mid n \geq 1\}$$

$\xrightarrow{\downarrow}$
no S remains in the resulting string

Thus

$$L(G) = \{a^n b^n \mid n \geq 1\}$$

Example

Find a CFG generating the language of even length palindromes.

i.e. $L = \{ww^R \mid w \in \Sigma^* = \{a,b\}^*\}$

$$G = (\{S\}, \{a,b\}, P, S),$$

P is defined as $S \rightarrow aSa \mid bSb \mid \epsilon$

Arguments:

- i. ϵ is a palindrome
- ii. a, b are palindromes
- iii. if x is a palindromes, then axa, bxb are palindromes.

Example:

Construct a CFG generating

$$L = \{wcw^T : w \in \{a,b\}^*\}$$

$$G = (\{S\}, \{a,b,c\}, P, S)$$

where P is defined as $S \rightarrow c \mid aSa \mid bSb$

Argument:

Any string in L is generated by recursion as follows:

- i. $c \in L$
- ii. if $x \in L$, then $wxw^T \in L$

Example

A CFG for the regular language corresponding to the regular expression 00^*11^* ,

$$L = \{ 0^i 1^j \mid i, j > 0 \}$$

idea:

The language is the concatenation of two languages

- language of strings of 0's
- language of strings of 1's

$$G = (\{S, C, D\}, \{0, 1\}, P, S)$$

where P is defined as $\{S \rightarrow CD, C \rightarrow 0C|0, D \rightarrow 1D|1\}$

Example:

CFG that generates sentences as composed of noun and verb phrases.

$$S \rightarrow NP \ VP$$

$$NP \rightarrow \text{the } N$$

$$VP \rightarrow V \ NP$$

$$V \rightarrow \text{sings/eats}$$

$$N \rightarrow \text{cat/song/bird}$$

"the bird sings the song" \rightarrow accepted

"the song eats the cat" \rightarrow also accepted

- All legal sentences are generated, not just those that are meaningful — content-free

Example:

Complement of $\{0^i1^j : i,j \geq 0\}$

Idea:

partition the set of strings into three failure cases

1. strings not of the right form: somewhere there is a 1 followed by a 0
2. only zeroes
3. only ones

$$G = (\{S, A, B, C, D\}, \{0, 1\}, P, S)$$

where P has the following production rules:

$$S \rightarrow A | B | C$$

$A \rightarrow D01D$ (produces all strings with a 0 and 1 out of order)

$D \rightarrow 0D | 1D | \epsilon$ (produces all strings)

$B \rightarrow 0B | 0$ (produces strings of 0's)

$C \rightarrow 1C | 1$ (produces strings of 1's)

e.g.

Test whether $w \in L(G)$ where

$$w = 101001$$

$$S \Rightarrow A \Rightarrow D10D \Rightarrow \epsilon 10D = 10D$$

$$\Rightarrow 101D$$

$$\Rightarrow 1010D$$

$$\Rightarrow 10100D$$

$$\Rightarrow 101001D$$

$$\Rightarrow 101001\epsilon$$

$$\Rightarrow 101001$$

$$= w$$

$$\therefore w \in L(G) = \{w | w \in T^* \text{ and } S^* \xrightarrow{*} w\}$$

Example:

Find a CFG G for all binary strings with an even number of 0's

$$S \rightarrow 1S \mid 0A0S \mid \epsilon$$

$$A \rightarrow 1A \mid \epsilon$$

produces strings of 1's.

How to decompose such strings?

1. 1st symbol starts with 1, followed by even no. of 0's
2. 1st symbol starts with 0, then go to the next 0, followed even number of 0's

e.g. 101000 $\in L(G)$ as

$$\begin{aligned} S &\Rightarrow 1S \Rightarrow 10A0S \Rightarrow 101A0S \Rightarrow 1010S = 1010S \\ &\Rightarrow 10100A0 \Rightarrow 10100\epsilon 0 \Rightarrow 101000 \end{aligned}$$

- A language can be more than one grammar

$S \rightarrow S \mid OT \mid \epsilon, T \rightarrow T \mid 0S \}$ also generates all binary strings with even no. of 0's

idea:

We decompose all binary strings with an even number of 0's as follows:

1. 1st symbol is 0, followed by odd number of
2. We use T to generate all binary strings with odd number of zeroes.
3. We use S to generate all binary string with even number of zeroes.

Production rules: $S \rightarrow 1S \mid OT \mid \epsilon$

$$T \rightarrow IT \mid 0S$$

Example:

Let $G = (\{S, C\}, \{a, b\}, P, S)$ where P consists of
 $S \rightarrow aCa, C \rightarrow aCa|b$. Find $L(G)$.

$$S \Rightarrow aCa \Rightarrow aba \quad \therefore aba \in L(G)$$

$$S \Rightarrow aCa \Rightarrow aaCaa \Rightarrow \dots \Rightarrow a^n Ca^n \Rightarrow a^n ba^n$$

$\therefore a^n ba^n \in L(G), n \geq 1$

$$\therefore \{a^n ba^n | n \geq 1\} \subseteq L(G) \quad \text{--- (1)}$$

$$\text{Now, to prove } L(G) \subseteq \{a^n b^n a^n | n \geq 1\}$$

Any string $w \in T^* = \{a, b\}^*$ will be in $L(G)$ if

$$S \xrightarrow{*} w$$

The only S -production is $S \rightarrow aCa$

if we apply $C \rightarrow b$, we get $aba \in T^*$

$$S \xrightarrow{*} aba$$

if we apply $C \rightarrow aCa$, once or several times

$$\text{we get } S \xrightarrow{*} a^n Ca^n$$

$$\Rightarrow a^n ba^n \in T^*$$

(one application of $C \rightarrow b$)

Thus any derivation is of the form

$$S \xrightarrow{*} a^n b a^n, n \geq 1$$

$$\Rightarrow L(G) \subseteq \{a^n ba^n | n \geq 1\} \quad \text{--- (2)}$$

Thus from (1) and (2)

$$L(G) = \{a^n ba^n | n \geq 1\}$$

Example:

Let $G = (\{S, A_1\}, \{0, 1, 2\}, P, S)$ where P consists of
 $S \rightarrow 0SA_12, S \rightarrow 012, 2A_1 \rightarrow A_12, 1A_1 \rightarrow 11$. Find $L(G)$

We have

$$S \Rightarrow 012 \quad \therefore 012 \in L(G)$$

$$S \Rightarrow 0\underline{S}A_12 \Rightarrow 001\underline{2}A_12 \Rightarrow 001\underline{A_1}22 \Rightarrow 001122$$

$$S \Rightarrow 0SA_12 \Rightarrow 00SA_12A_12 \Rightarrow 000SA_12A_12A_12$$

$S \Rightarrow 0SA_12$ applied $(n-1)$ times

$$\Rightarrow 0^{n-1} S (A_12)^{n-1}$$

$$\Rightarrow 0^{n-1} 012 (A_12)^{n-1}$$

$$\Rightarrow 0^n 1 2A_1 (A_1)^{n-2} 2^{n-1}$$

$$\Rightarrow 0^n 1 A_1 2A_1 (A_1)^{n-3} 2^{n-1}$$

$$\Rightarrow 0^n 1 A_1 A_1 2A_1 (A_1)^{n-4} 2^{n-1}$$

$$\vdots \Rightarrow 0^n 1 A_1^{n-1} 2^n$$

$$\Rightarrow 0^n 1 A_1 A_1^{n-2} 2^n$$

$$\Rightarrow 0^n 1 A_1 A_1 A_1^{n-3} 2^n$$

$$\vdots \Rightarrow 0^n 1^n 2^n$$

$$\therefore 0^n 1^n 2^n \in L(G)$$

$$\Rightarrow \{0^n 1^n 2^n \mid n \geq 1\} \subseteq L(G)$$

Then we need to show that $L(G) \subseteq \{0^n 1^n 2^n \mid n \geq 1\}$

Example

Let $G = (\{S, A_1, A_2\}, \{a, b\}, P, S)$

where P consists of

$S \rightarrow aA_1, A_2a, A_1 \rightarrow baA_1, A_2b$

$A_2 \rightarrow A_1ab, A_1 \rightarrow baa$

$bA_2b \rightarrow abab$

Test whether $w = baabbabaaaabbaba$ is in $L(G)$.

$S \Rightarrow a\underline{A_1} A_2a$

$\Rightarrow baa\underline{A_2}a$

$\Rightarrow baa\underline{A_1}aba$

$\Rightarrow baa ba\underline{A_1} A_2 baba$

$\Rightarrow baabbbaa\underline{A_2} baba$

$\Rightarrow baabbbaa\underline{aA_1}abbaba$

$\Rightarrow baabbabaaaabbaba = w$

$\therefore w \in L(G)$

Example

If the grammar G is given by the productions
 $S \rightarrow aSa \mid bSb \mid aa \mid bb \mid \epsilon$, show that

- i) $L(G)$ has no strings of odd length
- ii) any string in $L(G)$ is of length $2n$, $n > 0$
- iii) the number of strings of length $2n$ is 2^n

On application of each production (except $S \rightarrow \epsilon$), a variable is replaced by 2 terminals and atmost one variable

Thus each step in any derivation increases the number of terminals by 2 except that involving $S \rightarrow \epsilon$.

Which explains (i) and (ii)

Now any string w of length $2n$ is of the form

$$a_1 a_2 \dots a_n a_n a_{n-1} \dots a_1$$

where each a_i is either a or b

\Rightarrow number of $2n$ -bit strings in $L(G)$ is 2^n

Example

Construct a context free grammar G generating all integers (with sign)

Let $G = \{V, T, P, S\}$ where

$$V = \{S, \langle \text{sign} \rangle, \langle \text{digit} \rangle, \langle \text{integer} \rangle\}$$

$$T = \{0, 1, \dots, 9, +, -\}$$

and P consists of

$$S \rightarrow \langle \text{sign} \rangle \langle \text{integer} \rangle$$

$$\langle \text{sign} \rangle \rightarrow + | -$$

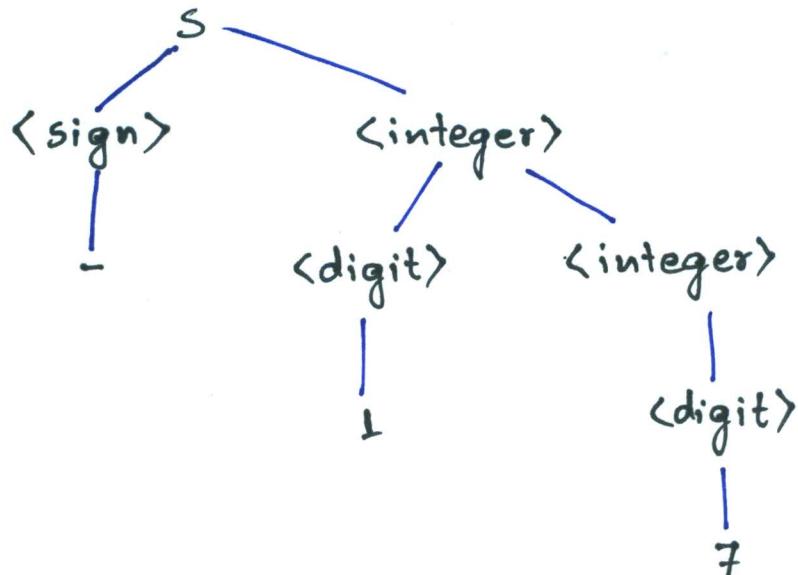
$$\langle \text{integer} \rangle \rightarrow \langle \text{digit} \rangle \langle \text{integer} \rangle | \langle \text{digit} \rangle$$

$$\langle \text{digit} \rangle \rightarrow 0 | 1 | 2 | \dots | 9$$

Then

$L(G) =$ the set of all integers.

e.g. derivation of -17



Derivation Trees / Parse Trees

- $G = (V, T, P, S)$ be a CFG. A tree is a derivation tree or parse tree for G if:
 - Every vertex has a label, which is a symbol of $V \cup T \cup \{\epsilon\}$
 - The label of the root is S
 - If a vertex is interior and has label A , then A must be in V
 - If vertex n has label A and vertices n_1, n_2, \dots, n_k are sons of vertex n , in order from the left with labels x_1, x_2, \dots, x_k respectively, then $A \rightarrow x_1 x_2 \dots x_k$ must be a production in P
 - If vertex n has label ϵ , then n is a leaf and is the only son of its father

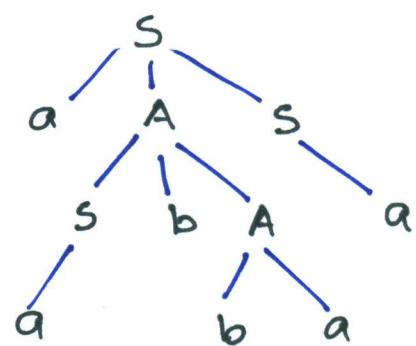
Example:

$$G = (\{S, A\}, \{a, b\}, P, S)$$

where consists of

$$S \rightarrow aAS \mid a$$

$$A \rightarrow SbA \mid \epsilon \mid ba$$



A natural description of
the sentential form

yield of the
derivation
tree $\rightarrow aabbbaaa$

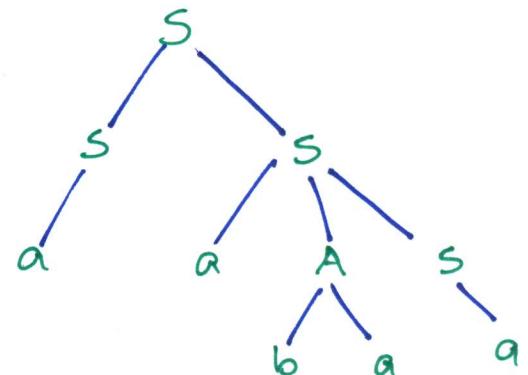
- The yield of a derivation tree is the concatenation of the labels of the leaves without repetition in the left-to-right ordering.

Example:

$$G = (\{S, A\}, \{a, b\}, P, S)$$

where P consists of

$$\begin{aligned} S &\rightarrow aAS \mid a \mid ss \\ A &\rightarrow SbA \mid ba \end{aligned}$$



A parse tree showing the derivation

$$S \xrightarrow{*} aabaa$$

- The yield of the derivation tree is a sentential form in G . i.e. $\text{yield} \in (VUT)^*$

Example:

$$G = (\{E, I\}, T, P, E)$$

where P consists of

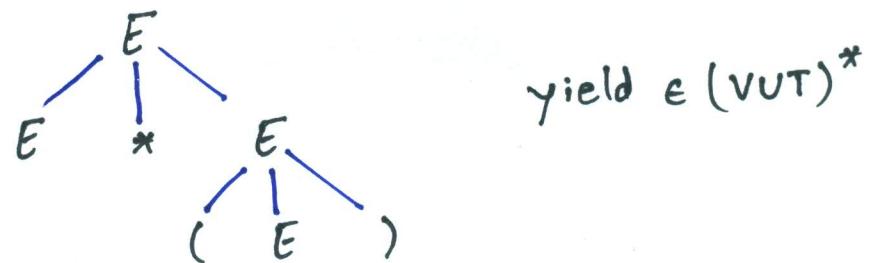
$$E \rightarrow I \mid E+E \mid E * E \mid (E)$$

$$I \rightarrow a \mid b \mid Ia \mid Ib \mid I_0 \mid I_1$$

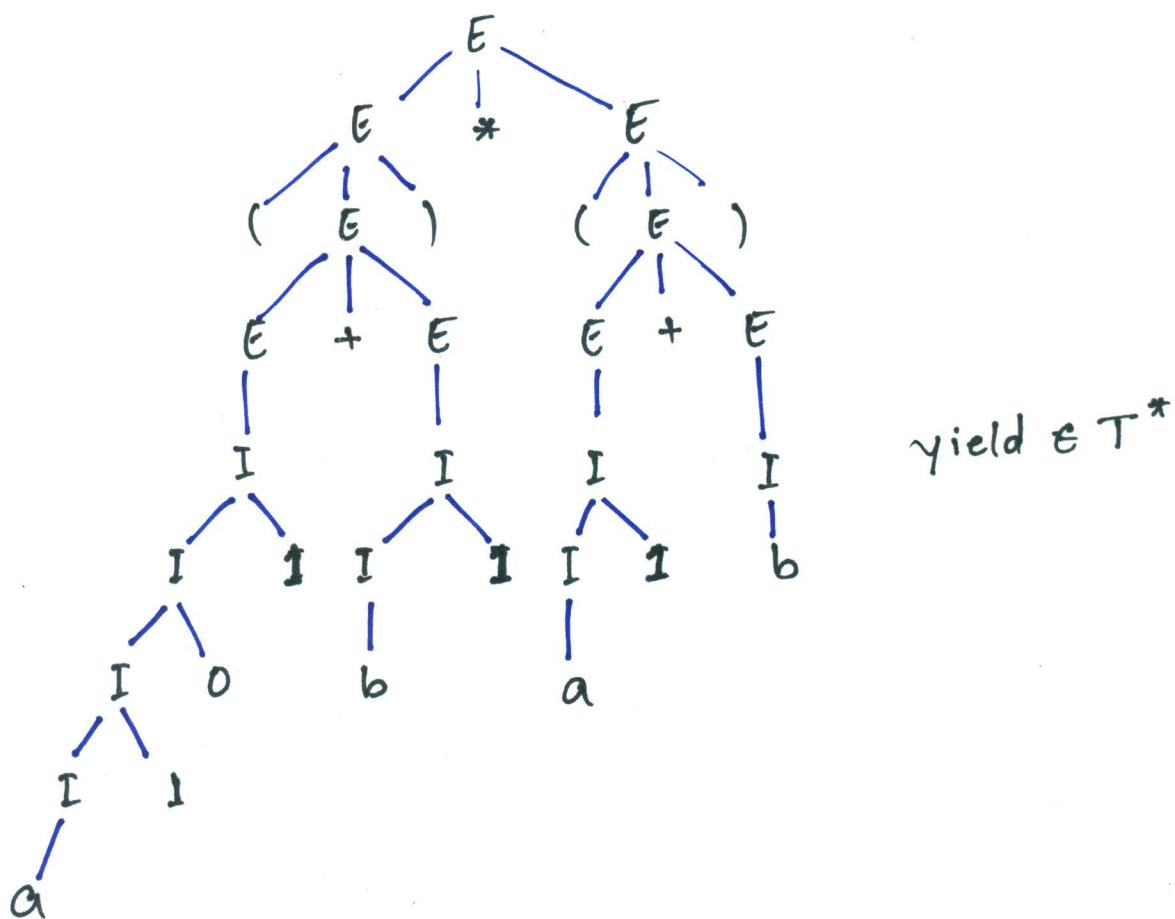
and

$$T = \{+, *, (,), a, b, 0, 1\}$$

(i) A parse tree showing the derivation
 $E \xrightarrow{*} E * (E)$



A parse tree showing the derivation
 $E \xrightarrow{*} (a_1 o_1 + b_1) * (a_1 + b)$



Relationship between derivation trees and derivations

Theorem:

Let $G = (V, T, P, S)$ be a context-free grammar. Then $S \xrightarrow{*} \alpha$ iff there is a derivation tree in grammar G with yield α

Proof:

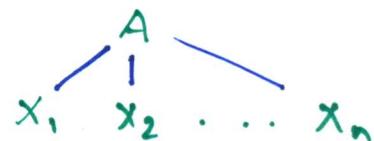
A-tree: A subtree of a derivation tree is called an A-tree if the label of its root is A

S-tree: derivation tree when S is the start symbol

We prove that $A \xrightarrow{*} \alpha$ iff there is an A-tree with yield α .

(if) Suppose α is the yield of an A-tree. We prove by induction on K = the no. of internal vertices in the tree, that $A \xrightarrow{*} \alpha$

Basis: $K=1$, i.e. only one internal vertex



$\therefore \alpha = x_1 x_2 \dots x_n$ and

$A \rightarrow \alpha$ must be a production rule in G by definition of a derivative tree

Induction step:

Let us assume the result for all derivation trees with $\leq k-1$ internal vertices for some $k \geq 1$

Let α be the yield of an A-tree with k -interior vertices for some $k \geq 1$

Let the sons of the root A have labels x_1, x_2, \dots, x_n in order from the left

Then $A \rightarrow x_1 x_2 \dots x_n$ is in P ($n \geq 1$ is any integer)

- x_1, x_2, \dots, x_n cannot all be leaves (as $k \geq 1$)
- if i -th son of A is not a leaf, then it is the root of a subtree and x_i must be a variable.

x_i -tree with yield α_i say.

if i -th son of A is a leaf, let $\alpha_i = x_i$

Then $\alpha = \alpha_1 \alpha_2 \dots \alpha_n$ as one can easily see that if $j < i$, then vertex j and all its descendants are to the left of vertex i and all its descendants.

Now subtrees must have fewer interior nodes than its tree does, unless the subtree is the entire tree.

∴ We have for each x_i -tree with yield α_i :

$x_i \xrightarrow{*} \alpha_i$ by induction as x_i -tree is not the entire tree.

Also, $x_i \Rightarrow \alpha_i$ if $x_i = \alpha_i$

$$\begin{aligned}\therefore A \Rightarrow x_1, x_2, \dots, x_n &\Rightarrow \alpha_1, \alpha_2, \dots, \alpha_n \\ &\stackrel{*}{\Rightarrow} \alpha_1, \alpha_2, \dots, x_n \Rightarrow \dots \\ &\stackrel{*}{\Rightarrow} \alpha_1, \alpha_2, \dots, \alpha_n = \alpha\end{aligned}$$

Thus $A \stackrel{*}{\Rightarrow} \alpha$

(Only if)

Let $A \stackrel{*}{\Rightarrow} \alpha$

We must show that there is an A-tree with yield α
We do it by induction on K : the number of steps
in $A \stackrel{*}{\Rightarrow} \alpha$

(Base)

$K=1$

If $A \stackrel{*}{\Rightarrow} \alpha$ in single step, then $A \rightarrow \alpha$ is in.

Induction:

Assume the result is true for derivatives in $< K$ steps

Let $A \stackrel{K}{\Rightarrow} \alpha$

This can be splitted as

$\underbrace{A \Rightarrow x_1, x_2, \dots, x_n}_{\text{implies}} \stackrel{K-1}{\Rightarrow} \alpha$

$A \rightarrow x_1, x_2, \dots, x_n$ is in P

In $x_1, x_2, \dots, x_n \stackrel{K-1}{\Rightarrow} \alpha$ either

i. x_i is not changed throughout the derivation

or ii. x_i is changed in some subsequent step.

Let α_i be the substring derived from x_i

Then $x_i = \alpha_i$ is in (i) and $x_i \xrightarrow{*} \alpha_i$ is in (ii)

- As G is content free, in every step of derivation

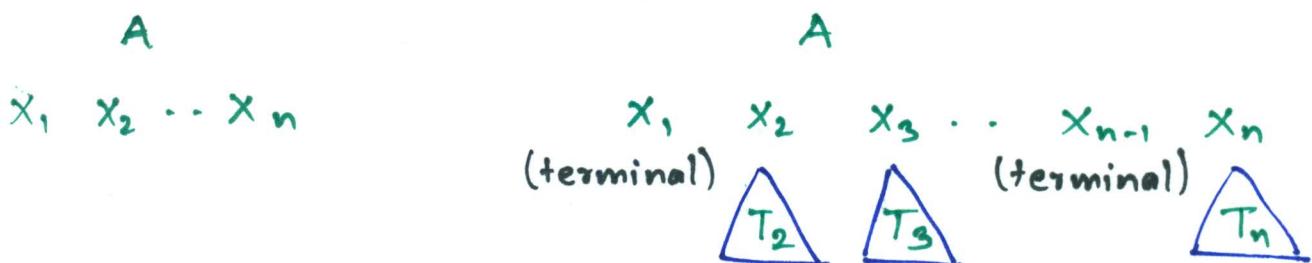
$x_1 x_2 \dots x_n \xrightarrow{*} \alpha$, we replace a single variable by a string

As $\alpha_1, \alpha_2, \dots, \alpha_n$ accounts for all the symbols in α
we have $\alpha = \alpha_1 \alpha_2 \dots \alpha_n$

- If x_i is a variable, then the derivation $x_i \xrightarrow{*} \alpha_i$ has $< K$ steps as the entire derivation $A \xrightarrow{*} \alpha$ takes K steps and the first step is surely not a part of $x_i \xrightarrow{*} \alpha_i$

Hence by induction, for each x_i that is a variable, there is an x_i -tree T_i with yield α_i

Now, we construct an A-tree with yield α as follows:



First construct an A-tree with n leaves labeled x_1, x_2, \dots, x_n and no other vertex X

Replace each vertex with label x_i by tree T_i where x_i is non-terminal.

If x_i is a terminal vertex, no replacement is made.

The yield of the resulting A-tree is $\alpha_1 \alpha_2 \dots \alpha_n = \alpha$

Note:

The derivation tree does not specify the order in which we apply the productions for getting α . So, the same derivation tree can induce several derivations (with the same yield)

Remark:

Let $A \xrightarrow{*} w$ where $w \in T^*$ (a terminal string)
and $A \Rightarrow A_1 A_2 \dots A_n$

Then we can write $w = w_1 w_2 \dots w_n$
so that $A_i \xrightarrow{*} w_i$

Week 9: Lecture Notes.

Theorem:

Let $G = (V, T, P, S)$ be a context free grammar. Then $S \xrightarrow{*} \alpha$ iff there is a derivation tree in grammar G with yield α .

Example:

Consider CFG G whose productions are

$S \rightarrow aAS | a, A \rightarrow SBa | ss | ba$. Show that $S \xrightarrow{*} aabbbaa$ and construct a derivation tree whose yield is $aabbbaa$

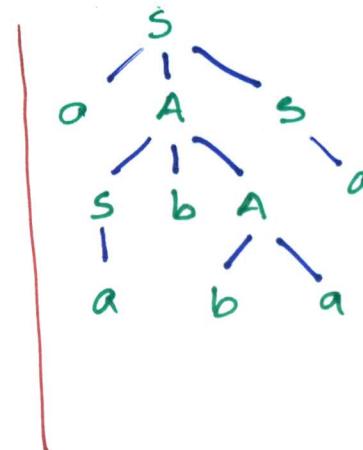
Sol:

$$S \xrightarrow{*} aAS \Rightarrow a \underline{S} bAS \Rightarrow aab \underline{A} S \Rightarrow aabb \underline{a} \xrightarrow{*} aabbbaa$$

another derivation:

$$S \xrightarrow{*} aAS \Rightarrow aAa$$

$$\begin{aligned} &\Rightarrow aSbAa \Rightarrow aabAa \Rightarrow aabbbaa \\ \Rightarrow &aSbAa \Rightarrow aSbbbaa \Rightarrow aabbbaa \end{aligned}$$



Leftmost and rightmost derivations: ambiguity

- A derivation $A \xrightarrow{*} w$ is called a leftmost derivation if we apply a production only to the leftmost variable at every step
- A derivation $A \xrightarrow{*} w$ is a rightmost derivation if we apply a production only to the rightmost variable at every step.

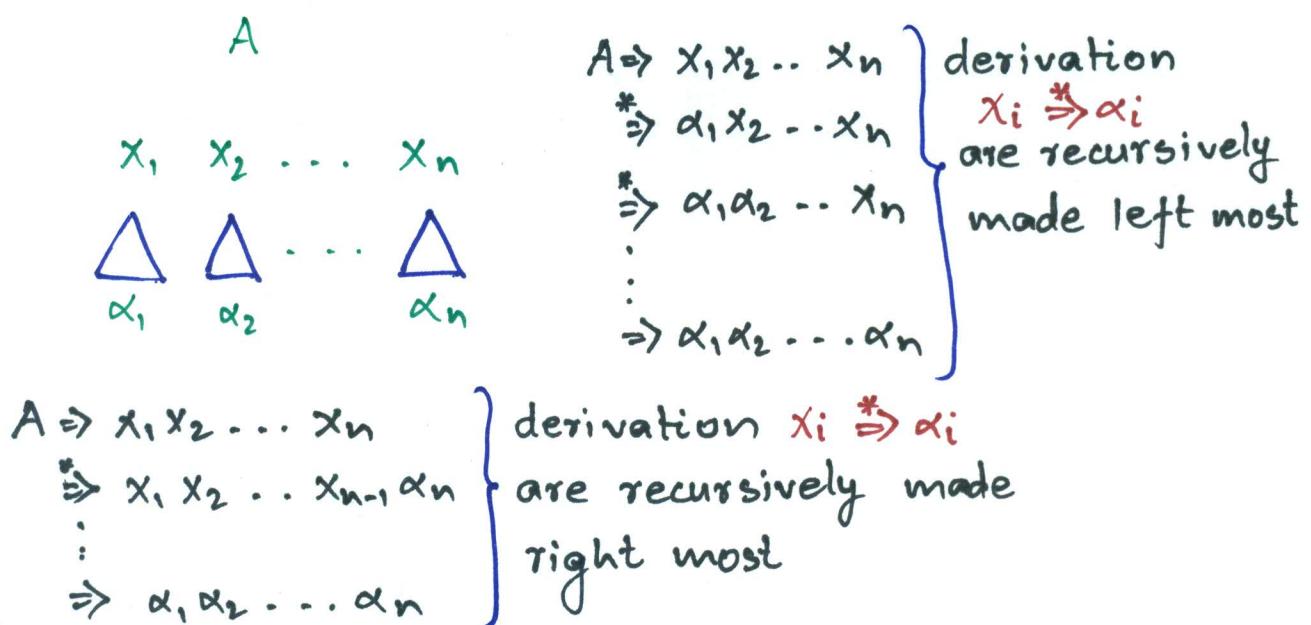
Theorem:

If $A \xrightarrow{*} w$ in G , then there is a leftmost derivation of w

Proof:

By induction on the number of steps in $A \xrightarrow{*} w$, we can prove it.

- If $w \in L(G)$ for CFG G , then w has atleast one parse tree and corresponding to a particular parse tree, w has a unique leftmost and a unique rightmost derivation



- w may have several leftmost or rightmost derivations as they may be more than one parse tree for w

Example:

Let G be the CFG. $S \rightarrow 0B|IA$, $A \rightarrow 010S|1AA$, $B \rightarrow 1|1S|0BB$.

for the string 00110101 , find

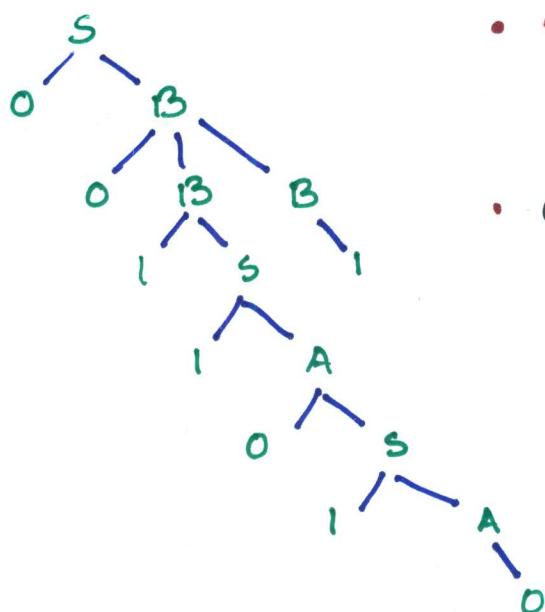
- the leftmost derivation
- the rightmost derivation
- the derivation tree

a). $S \Rightarrow 0B \Rightarrow 00BB \Rightarrow 001sB \Rightarrow 0011AB \Rightarrow 001101A \Rightarrow$
 $\Rightarrow 0011010B \Rightarrow 00110101$

b). Corresponding rightmost derivation

$S \Rightarrow 0B \Rightarrow 00B \Rightarrow 00B1 \Rightarrow 001s1 \Rightarrow 0011A1 \Rightarrow$
 $\Rightarrow 00110s1 \Rightarrow 001101A1 \Rightarrow 00110101$

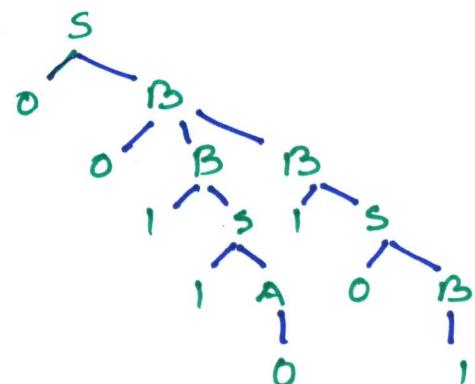
c). The derivation tree for w



- w may have several derivation trees
- Corresponding to each parse tree of w we have a unique leftmost derivation and a unique rightmost derivation

Another leftmost derivation:

$S \Rightarrow 0B \Rightarrow 00BB \Rightarrow 001sB$
 \Downarrow
 $001101s \Leftarrow 00110B \Leftarrow 0011AB$
 \Downarrow
 $0011010B \Rightarrow 00110101$

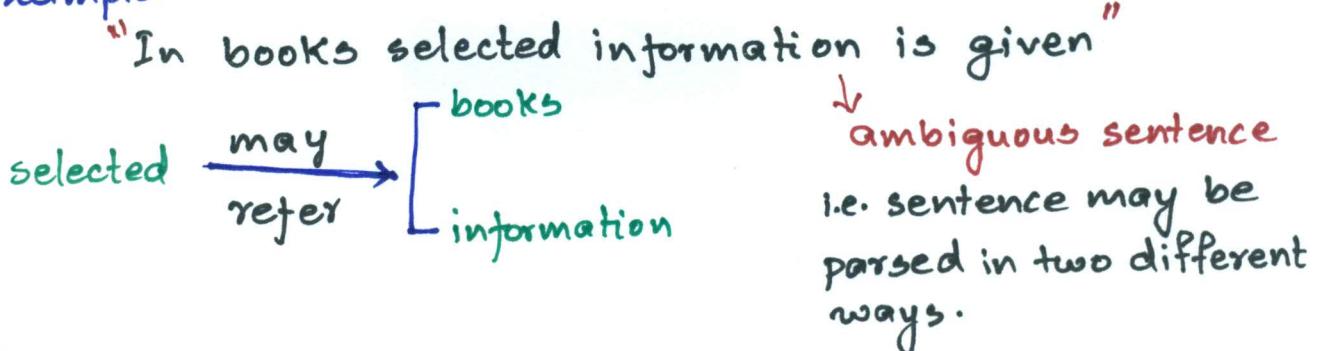


Corresponding rightmost derivation

$S \Rightarrow 0B \Rightarrow 00BB \Rightarrow 00B1s \Rightarrow 00B10B \Rightarrow 00B101$
 \Downarrow
 $00110101 \Leftarrow 0011A101 \Leftarrow 001s101$

Ambiguity in CFG

Example:



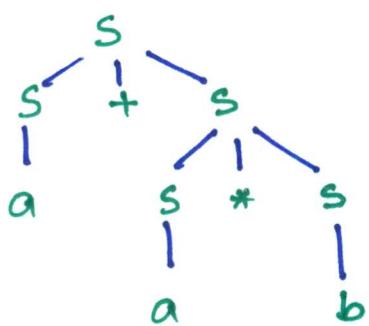
- A terminal string $w \in L(G)$ is ambiguous if \exists two or more derivation trees for w .
- A CFG G is ambiguous if \exists some $w \in L(G)$, which is ambiguous.
- A CFL for which every CFG is ambiguous is said to be an inherently ambiguous CFL.

Example:

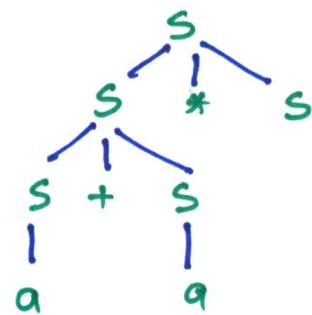
$G = (\{S\}, \{a, b, +, *\}, P, S)$, where P consists of

$$S \rightarrow S+S \mid S*S \mid a \mid b$$

Parse trees for $a+a*b$



T_1



T_2

The leftmost derivations of $a+a*b$ induced by T_1, T_2

- $S \Rightarrow S+S \Rightarrow a+S \Rightarrow a+S*S \Rightarrow a+a*S \Rightarrow a+a*b$
- $S \Rightarrow S*S \Rightarrow S+S*S \Rightarrow a+S*S \Rightarrow a+a*S \Rightarrow a+a*b$

Thus $a+a*b$ is ambiguous

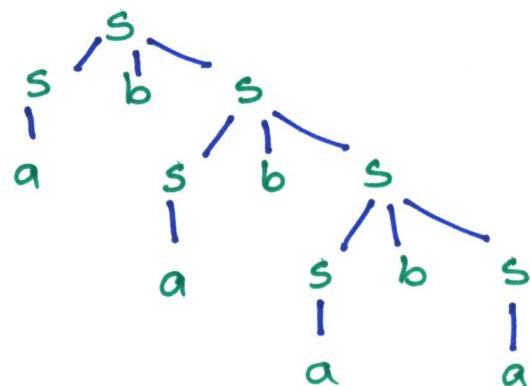
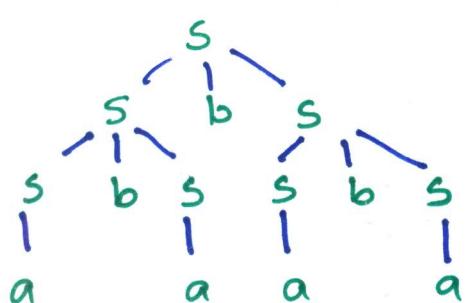
Example:

If G is the CFG $S \rightarrow SbS|a$. Show that G is ambiguous.

Solⁿ:

We have to find a $w \in L(G)$ which is ambiguous.

Let $w = abababa$



Simplification of CFG's

- $G = (V, T, P, S) \rightarrow \text{CFG}$
- $L(G)$ generation \rightarrow

may be not all symbols in $V \cup T$ are used

OR

may be not all rules in P are used
- Try to eliminate those symbols in $V \cup T$ or rules in P from G and construct a CFG G' that realizes $L(G)$

Example:

$$G = (\{S, A, B, D, E\}, \{a, b, c\}, P, S)$$

where $P = \{S \rightarrow AB, A \rightarrow a, B \rightarrow b, B \rightarrow D, E \rightarrow C|\epsilon\}$

Solⁿ:

Note that $L(G) = \{ab\}$

Let $G' = (\{S, A, B\}, \{a, b\}, P', S)$

where $P' = \{S \rightarrow AB, A \rightarrow a, B \rightarrow b\}$

WHY?

• eliminate symbols D, E, c • eliminate productions $B \rightarrow D, E \rightarrow C|\epsilon$

1. D does not derive any terminal string
2. E and c does not appear in any sentential form
3. $E \rightarrow \epsilon$ is a null production
4. $B \rightarrow D$ simply replaces B by D (unit production)

Algorithms to construct reduced grammars.

Construction 1

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

Define $G' = (V', T, P', S)$ as follows:

Construction of V'

We define $W_i \subseteq V$ by recursion:

$$W_i = \{ X \in V \mid \text{there exists a production } X \rightarrow w \text{ where } w \in T^* \}$$

[if $W_i = \emptyset$; then some variable will remain after application of any production and so $L(G) = \emptyset$]

$$W_{i+1} = W_i \cup \{ X \in V \mid \text{there exists a production } A \rightarrow \alpha \text{ with } \alpha \in (T \cup W_i)^* \}$$

Thus $W_i \subseteq W_{i+1} \forall i$

Also, $W_k = W_{k+j}$ for some $k \leq |V|$ as V is a finite set of variables.

Therefore $W_k = W_{k+j}$ for $j > 1$

We define $V' = W_k$

Construction of P'

$$P' = \{ A \rightarrow \alpha \mid A, \alpha \in (V' \cup T)^* \}$$

- $S \in V'$ (Prove it !!.)

Example:

Let $G = (V, T, P, S)$ be given by the productions
 $S \rightarrow AB$, $A \rightarrow a$, $B \rightarrow b$, $B \rightarrow C$, $E \rightarrow d$. find G' such that
every variable in G' derives some terminal string.

Sol:-

Construction of V'

$$W_1 = \{A, B, E\}$$

$$W_2 = W_1 \cup \{x \in V \mid \exists \text{ some production } X \rightarrow \alpha \text{ with } \alpha \in (W_1 \cup T)^*\}$$

$$= W_1 \cup \{S\}$$

$$= \{A, B, E, S\}$$

$$W_3 = W_2 \cup \{x \in V \mid \exists \text{ some production } X \rightarrow d \text{ with } \alpha \in (W_2 \cup T)^*\}$$

$$= W_2 \cup \emptyset$$

$$= W_2$$

$$\therefore V' = \{A, B, E, S\}$$

Construction of P'

$$P' = \{X \rightarrow \alpha \mid X, \alpha \in (V' \cup T)^*\}$$

$$= \{S \rightarrow AB, A \rightarrow a, B \rightarrow b, E \rightarrow d\}$$

$$\therefore G' = (\{S, A, B, E\}, \{a, b, d\}, \underline{\underline{P'}}, S)$$

Note:-

$S \in V'$ always

Theorem 1

Given a CFG $G = (V, T, P, S)$, with $L(G) \neq \emptyset$, we can effectively find an equivalent CFG $G' = (V', T, P', S)$ such that for each A in V' there is some w in T^* for which $A \xrightarrow[G']{*} w$

Proof:

Let V', P' are constructed from G following the above algorithm.

Claim:

- i) for each $A \in V'$, $A \xrightarrow[G']{*} w$ for some $w \in T^*$
conversely $A \xrightarrow[G']{*} w$ implies $A \in V'$
- ii) $L(G) = L(G')$

Proof of claim (i)

Note that $w_k = w_1 w_2 w_3 \dots w_k$

We prove by induction on i for $i = 1, 2, \dots, k$ that $A \in W_i$ implies $A \xrightarrow[G']{*} w$ for some $w \in T^*$

Base:

Let $A \in W_1$,

Then $A \rightarrow w$ is in P for some $w \in T^*$

By construction of P' , $A \rightarrow w$ is in P'

This in turn implies $A \xrightarrow[G']{*} w$

Induction:

Let the result be true for i

Consider $A \in W_{i+1}$

Then, either

- $A \in W_i$, in which case $A \xrightarrow[G]{*} w$ for some $w \in T^*$
by induction
- $\exists \alpha$ production $A \rightarrow \alpha$ in P with $\alpha \in (TUW_i)^*$
in which case $A \rightarrow \alpha$ is in P' by construction
of P'

Then α may be written as

$$\alpha = x_1 x_2 \dots x_n \text{ where } x_j \in TUW_i$$

If $x_j \in W_i$, then by induction hypothesis,

$$x_j \xrightarrow[G]{*} w_j \text{ for some } w_j \in T^*$$

Hence, $A \xrightarrow[G]{*} w_1 w_2 \dots w_n \in T^*$

Conversely we can show that $A \xrightarrow[G]{*} w$ implies $A \in V'$

We can prove it by induction on $K = \text{number of steps in the derivation } A \xrightarrow[G]{*} w$ in a similar manner.

Base: $K=1$

production $A \rightarrow w$ is in P , so $A \in V'$

Induction

Let $A \xrightarrow[\alpha]^* x_1 x_2 \dots x_n \xrightarrow[\alpha]^{K-1} w$ by a derivation of K steps.

We may write $w = w_1 w_2 \dots w_n$ where $x_i \xrightarrow{*} w_i$ for $1 \leq i \leq n$, by a derivation of fewer than K steps.

∴ By induction hypothesis $x_i \in V'$, $1 \leq i \leq n$

Hence, $A \rightarrow x_1 x_2 \dots x_n$ is in P' by construction of P' .

Let $V' = w_1 w_2 \dots w_l$ and $w_k = w_{k+1}$

Let j be the least integer, $1 \leq j \leq l$ for which w_j contains all x_1, x_2, \dots, x_n

As $A \rightarrow x_1 x_2 \dots x_n$ is in P , where $x_1, x_2, \dots, x_n \in (TUV_j)^*$ we must have $A \in W_{j+1} \subseteq V'$

Thus $A \xrightarrow[\alpha]^K w$ implies $A \in V'$

Proof of claim (ii)

• $L(G') \subseteq L(G)$ as $P' \subseteq P$ and $V' \subseteq V$ —— (1)

• To prove $L(G) \subseteq L(G')$, we need to prove the following auxiliary result.

$A \xrightarrow[G']{*} w$ if $A \xrightarrow[\alpha]{*} w$ for some $w \in T^*$

We prove this by induction on K = the number of steps in the derivation $A \xrightarrow[\alpha]{*} w$

Base:

$K=1$

Let $A \xrightarrow[G]{*} w$ where $w \in T^*$

Then $A \xrightarrow{P'} w$ is in P'

$\therefore A \xrightarrow[G']{*} w$

Induction

Result is true for $\leq K$ steps.

Let $A \xrightarrow[G]{K+1} w$ where $w \in T^*$. We need to prove $A \xrightarrow[G']{K+1} w$

We may write $w = w_1 w_2 \dots w_n$ and

$A \xrightarrow[G]{*} x_1 x_2 \dots x_n \xrightarrow[G]{K} w$ s.t. $x_j \xrightarrow[G]{*} w_j$, $1 \leq j \leq n$
is in P

- If $x_j \in T$ then $x_j = w_j$
- If $x_j \in V$ then $x_j \xrightarrow[G]{*} w_j$ implies $x_j \in V'$
by claim(i)

Also $x_j \xrightarrow[G]{*} w_j$ in at most K steps.

Hence by induction hypothesis $x_j \xrightarrow[G]{*} w_j$

Also $x_1 x_2 \dots x_n \in (V' \cup T)^*$ implies

$A \xrightarrow{P'} x_1 x_2 \dots x_n$ is in P'

Thus $A \xrightarrow[G]{*} x_1 x_2 \dots x_n \xrightarrow[G]{K} w_1 w_2 \dots w_n$

i.e. $A \xrightarrow[G]{K+1} w$

In particular, $S \xrightarrow{*} w$ for $w \in T^*$

This implies $L(G) \subseteq L(G')$ ————— (ii)

From (i) and (ii), we get $L(G) = L(G')$

Theorem 2.

Given a CFG $G = (V, T, P, S)$, we can effectively find an equivalent CFG $G' = (V', T', P', S)$ such that for each $x \in V'UT'$, there exists $\alpha, \beta \in (V'UT')^*$ for which

$$S \xrightarrow{*} \alpha x \beta$$

Proof:

The set $V'UT'$ of symbols appearing in sentential form of G is constructed by an iterative algorithm.

- Algorithm 2.
- Construction 2.
- a. Place S in V'
 - b. If A is placed in V' and $A \rightarrow \alpha_1|\alpha_2|...|\alpha_n$, then add all variables of $\alpha_1, \alpha_2, \dots, \alpha_n$ to the set V' and all terminals of $\alpha_1, \alpha_2, \dots, \alpha_n$ to T'
 - c. P' is the set of productions of P containing only symbols of $V'UT'$

Example:

Consider $G = (\{S, A, B, E\}, \{a, b, c\}, P, S)$ where P consists of $S \rightarrow AB, A \rightarrow a, B \rightarrow b, E \rightarrow c$. Construct an equivalent grammar $G' = (V', T', P', S)$ such that every symbol in $V'UT'$ appears in some sentential form of G' .

Soln:

0. Initially set $V' = \{S\}$
 $S \rightarrow AB$
1. $V' = \{S, A, B\}$.
 $S \rightarrow AB, A \rightarrow a, B \rightarrow b$
2. $V' = \{S, A, B\}, T' = \{a, b\}$

No further modification is possible as $S \rightarrow AB$, $A \rightarrow a$, $B \rightarrow b$.

$$\therefore V' = \{S, A, B\}, T' = \{a, b\}$$

$$P' = \{S \rightarrow AB, A \rightarrow a, B \rightarrow b\}$$

Algorithm 1:

input: CFG G with $L(G) \neq \emptyset$

output: equivalent CFG G' such that every variable in G' derives some terminal string

Algorithm 2:

input: CFG $G = (V, T, P, S)$

output: equivalent CFG $G' = (V', T', P', S)$ such that
every symbol in $V'UT'$ appears in some sentential form.

Definitions:

- G is said to be **reduced** or **non-redundant** if every symbol in VUT appears in the course of the derivation of some terminal string
i.e. for every $x \in VUT$, \exists some $\alpha, \beta \in (VUT)^*$ such that $S \xrightarrow{*} \alpha x \beta \xrightarrow{*} w$, where $w \in T^*$
- We say X is **useful** in the derivation of terminal string.

Theorem 3.

For every CFG G , there exists a reduced grammar G' which is equivalent to G

Proof:

We construct the reduced grammar in two steps:

- | | |
|--|---|
| <p>Algorithm 3</p> <p>Construction 3</p> | <p>Step 1:
Construct a CFG G_1 equivalent to G using Algorithm 1, so that every variable in G_1 derives some terminal string</p> <p>Step 2:
Construct a CFG G' equivalent to G_1 using Algm.2 so that every symbol (variable as well as terminal symbol) in G' appears in some sentential form of G'</p> |
|--|---|

By step 2, every symbol X in G' appears in some sentential form, say $\alpha X \beta$

By step 1, every symbol in $\alpha X \beta$ derives some terminal string.

Note that all symbols of G' are symbols of G

It follows that $S \xrightarrow{*} \alpha X \beta \xrightarrow{*} w$ for some terminal string w

Hence G' is reduced.

NOTE:

To get a reduced grammar, we must first apply Algorithm 1 and then Algorithm 2. For if Algorithm 2 is applied first and then Algorithm 1, we may not get a grammar with its most reduced form

Example:

Find a reduced grammar equivalent to the grammar G whose productions are

$$S \rightarrow AB \mid CA, B \rightarrow BC \mid AB, A \rightarrow a, C \rightarrow aB \mid b$$

Solⁿ:

$$G = (V, T, P, S) \xrightarrow[1]{Alg'm} G_1 = (V', T', P', S) \xrightarrow[2]{Alg'm} G_2 = (V'', T'', P'', S)$$

G_1

$W_1 = \{A, c\}$ as $A \rightarrow a, C \rightarrow b$ are productions with a terminal string on right hand side

$W_2 = \{A, c\} \cup S$ as $S \rightarrow CA$ is a production with R.H.S. $\in \{T \cup W_1\}^*$

$W_3 = \{A, c, S\} \cup \emptyset$ as $\{x \in V \mid x \rightarrow \alpha \text{ with } \alpha \in (T \cup W_2)^*\}$ is empty

$$\therefore V' = \{A, C, S\}, P' = \{A \rightarrow a, C \rightarrow b, S \rightarrow CA\}$$

G_2

0. $V'' = \{S\}$

$$S \rightarrow CA$$

1. $V'' = \{S, C, A\}$

$$S \rightarrow CA, C \rightarrow b, A \rightarrow a$$

2. $V'' = \{S, C, A\}, T'' = \{a, b\}$

$$P'' = \{S \rightarrow CA, C \rightarrow b, A \rightarrow a\}$$

\therefore The reduced equivalent grammar is

$$G' = (\{S, A, C\}, \{a, b\}, \{S \rightarrow CA, C \rightarrow b, A \rightarrow a\}).$$

Example:

Construct a reduced grammar with no useless symbol equivalent to the grammar

$$S \rightarrow aAa, A \rightarrow Sb | bcc | Da, C \rightarrow abb | DD, E \rightarrow aC, D \rightarrow aDA$$

Sol:

Step 1: $W_1 = \{c\}$ as $C \rightarrow abb$ is the only production with a terminal string on the R.H.S.

$W_2 = \{c\} \cup \{A, E\}$ as $A \rightarrow bcc$ and $E \rightarrow aC$ are productions with R.H.S. in $(TUW_1)^*$

$W_3 = W_2 \cup \{S\}$ as $S \rightarrow aAa$ and $aAa \in (TUW_2)^*$

$W_4 = W_3 \cup \emptyset$

$$\therefore V' = W_4 = \{A, E, C, S\}$$

$$P' = \{A_i \rightarrow \alpha \mid \alpha \in (V' \cup T)^*\}$$

$$= \{S \rightarrow aAa, A \rightarrow Sb | bcc, C \rightarrow abb, E \rightarrow aC\}$$

Step 2:

0. $V'' = \{S\}$

$$S \rightarrow aAa$$

1. $V'' = \{S, A\}, T'' = \{a\}$

$$S \rightarrow aAa, A \rightarrow Sb | bcc$$

2. $V'' = \{S, A, C\}, T'' = \{a, b\}$

$$S \rightarrow aAa, A \rightarrow Sb | bcc, C \rightarrow abb$$

3. $V'' = \{S, A, C\}, T'' = \{a, b\}$

$$P'' = \{S \rightarrow aAa, A \rightarrow Sb | bcc, C \rightarrow abb\}$$

Algorithm 1: Transforms G to an equivalent G' where every variable derives some terminal strings.

Algorithm 2: Transforms G to an equivalent G' where every symbol appears in some sentential form.

Algorithm 3: Transforms G to an equivalent G' where G' is reduced, i.e. have no useless symbols.

Example:

$$S \rightarrow aAa, A \rightarrow Sb | bcc | DaA, C \rightarrow abb | DD, E \rightarrow aC \\ D \rightarrow aDA$$

Algorithm 1:

$$G': V' = W_0 = \{c\} \\ W_1 = \{c\} \cup \{A, E\} = \{c, A, E\} \\ W_2 = \{A, c, E\} \cup \{S\} \\ = \{S, A, C, E\} = V'$$

$$P': S \rightarrow aAa, A \rightarrow Sb | bcc, C \rightarrow abb, E \rightarrow aC$$

Algorithm 2:

$$G'': V'' = \{S\}, S \rightarrow aAa \\ V'' = \{S, A\}, T'' = \{a\}, A \rightarrow Sb | bcc \\ V'' = \{S, A, C\}, T'' = \{a, b\}, C \rightarrow abb$$

$$V'' = \{S, A, C\}, T'' = \{a, b\},$$

$$P'' = \{S \rightarrow aAa, A \rightarrow Sb | bcc, C \rightarrow abb\}$$

Algorithm 4: Transforms G to G' with no null production
where $L(G') = L(G) - \{\epsilon\}$

Algorithm 5: Transforms G to an equivalent G' with no null productions or unit productions

Algorithm 6: Transforms G to an equivalent reduced G' with no null productions or unit productions

Algorithm 4

$$G \quad G_1 = (V_1, T_1, P_1, S_1)$$

$$L(G) = L(G) - \{\epsilon\}$$

$$\downarrow \\ G_2 \text{ (no null production)}$$

$$L(G_2) = L(G)$$

$$\downarrow \\ G_3 \text{ (no unit production)}$$

$$\downarrow \\ \text{Algorithm 1}$$

$$G_4 \text{ (all variables derive some terminal string)}$$

$$\downarrow \\ \text{Algorithm 2}$$

$$G_5 \text{ (every symbol appears in some sentential form)}$$

- if $\epsilon \notin L(G)$, then $L(G_1) = L(G)$ set $G_2 = G_1$,

- if $\epsilon \in L(G)$ then $L(G_2) = L(G)$ where

$$G_2 = (V_2, T_2, P_2, S_2), \quad V_2 = V \cup \{S_2\}$$

$$P_2 = P_1 \cup \{S_2 \rightarrow \epsilon | S_1\}$$

Reduction to CNF

1. Eliminate null productions and unit productions
2. Elimination of terminals on R.H.S. of each productions
 - i. $A \rightarrow a$ $A \rightarrow BC$ type productions are retained
 - ii. $A \rightarrow X_1 X_2 \dots X_n$ type productions:
if some $X_i = a_i$, a terminal, introduce a new variable C_{a_i} and introduce a production $C_{a_i} \rightarrow a_i$
3. Restricting the number of variables on R.H.S.

$A \rightarrow A_1 A_2 \dots A_n$, $n > 3$ all A_i variables.

$A \rightarrow A, C_1, C_1 \rightarrow A_2 C_2, \dots, C_{n-1} \rightarrow A_{n-1} A_n$

Elimination of Null productions (ϵ -productions)

- A production of the form $A \rightarrow \epsilon$ is called a **null production** where A is a variable

A variable A in a CFG is **nullable** if

$$A \xrightarrow{*} \epsilon$$

Theorem 4

If $L = L(G)$ ($\epsilon \notin L$) for some CFG $G = (V, T, P, S)$ then $L - \{\epsilon\}$ is $L(G')$ for a CFG G' with no null productions

Proof:

We construct $G' = (V, T, P', S)$ as follows:

Construction 4. We find the nullable variables recursively:
(i) $W_i = \{A \in V \mid A \rightarrow \epsilon \text{ is in } P\}$
(ii) $W_{i+1} = W_i \cup \{A \in V \mid \text{there exists a production } A \rightarrow \alpha \text{ with } \alpha \in W_i^*\}$

Step 1: Construction of the set of nullable variables

By definition, $W_i \subseteq W_{i+1} \ \forall i$

As V is finite, $W_{k+1} = W_k$ for some $k < |V|$

$\therefore W_{k+j} = W_k \ \forall j$

Let $W = W_k$

W is set of all nullable variables

i.e. for every variable X in W we have

$$X \xrightarrow[G]{*} \epsilon$$

Step 2: Construction of P'

- Any production whose R.H.S. does not have any nullable variable is included in P'
- If $A \rightarrow x_1 x_2 \dots x_n$ is in P , then add all productions $A \rightarrow \alpha_1 \alpha_2 \dots \alpha_n$ in P' where
 - i. $\alpha_i = x_i$ if $x_i \notin W$
 - ii. $\alpha_i = x_i \text{ or } \epsilon$ if $x_i \in W$
 - iii. $\alpha_1 \alpha_2 \dots \alpha_n = \epsilon$ i.e. not all α_i 's are ϵ

Thus productions of P' are obtained by

- either not erasing any nullable variables on the R.H.S. of $A \rightarrow x_1 x_2 \dots x_n$ (in P)
- or erasing some or all nullable variables provided some symbol appear on the R.H.S. of $A \rightarrow x_1 x_2 \dots x_n$ after erasing.

Let $G' = (V, T, P', S)$

Claim: G' has no null productions

i.e. $\nexists A \in V, w \in T^*$

$A \xrightarrow[G']{*} w$ iff $w \neq \epsilon$ and $A \xrightarrow[G]{*} w$

Example:

Consider the grammar G whose productions are

$$S \rightarrow aS \mid AB, A \rightarrow \epsilon, B \rightarrow \epsilon, D \rightarrow b$$

Construct a grammar G' without null productions generating $L(G) - \{\epsilon\}$.

Sol:

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

$$W_1 = \{A, B\}$$

$$W_2 = W_1 \cup \{A, GV \mid A_i \rightarrow \alpha \text{ is in } P \text{ with } \alpha \in W_1^*\}$$

$= \{A, B, S\}$ as $S \rightarrow AB$ is in P with $AB \in W_1^*$

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

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

Step 2: Construction of P'

1. $D \rightarrow b$ is included in P'

2. $S \rightarrow aS$ gives rise to $S \rightarrow aS$ and $S \rightarrow a$

3. $S \rightarrow AB$ gives rise to $S \rightarrow AB, S \rightarrow A, S \rightarrow B$

\therefore The required CFG without null productions is

$$G' = (\{S, A, B\}, \{a, b\}, P', S)$$

where P' consists of

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

$$S \rightarrow A, S \rightarrow B$$

To prove $L(G') = L(G) - \{\epsilon\}$ we prove an auxiliary result given by the following relation

$\forall A \in V$ and $w \in T^*$

$$A \xrightarrow[G']{*} w \text{ iff } w \neq \epsilon \text{ and } A \xrightarrow[G]{*} w$$

Proof:

(if) Let $A \xrightarrow[G']{*} w$ and $w \neq \epsilon$

We prove that $A \xrightarrow[G']{*} w$ by induction on the number

of steps in the derivation of $A \xrightarrow[G']{*} w$

Basis:

If $A \xrightarrow[G]{*} w$ and $w \neq \epsilon$ then $A \rightarrow w$ is in P' and

so $A \xrightarrow[G']{*} w$

Induction:

Assume the result for derivation in at most i steps.

Let $A \xrightarrow[G]{i+1} w$ and $w \neq \epsilon$

We can split the derivation as

$$A \xrightarrow[G]{} x_1 x_2 \dots x_n \xrightarrow[G]{i} w_1 w_2 \dots w_n$$

where

$$w = w_1 w_2 \dots w_n$$

and $x_j \xrightarrow[G]{*} w_j$

As $w \neq \epsilon$, not all w_j 's are ϵ

- if $w_j \neq \epsilon$, then by induction hypothesis
 $x_j \xrightarrow[\mathcal{G}^*]{} w_j$ and $w_j \neq \epsilon$ implies $x_j \xrightarrow[\mathcal{G}'^*]{} w_j$
- if $w_j = \epsilon$ then $x_j \in W$ (i.e. x_j is a nullable variable)
so using production $A \rightarrow x_1 x_2 \dots x_n$ in P
we construct production $A \rightarrow \alpha_1 \alpha_2 \dots \alpha_n$ in P'
where

$$\alpha_j = \begin{cases} x_j & \text{if } w_j \neq \epsilon \\ \epsilon & \text{if } w_j = \epsilon \text{ (i.e. } x_j \in W\text{)} \end{cases}$$

Since $w \neq \epsilon$ not all α_j 's are ϵ

$$\begin{aligned} \therefore A \xrightarrow[\mathcal{G}^*]{} \alpha_1 \alpha_2 \dots \alpha_n &\xrightarrow[\mathcal{G}'^*]{} w_1 w_2 \dots w_n \\ &\xrightarrow[\mathcal{G}^*]{} w_1 w_2 \alpha_3 \dots \alpha_n \\ &\vdots \\ &\xrightarrow[\mathcal{G}^*]{} w_1 w_2 \dots w_n = w. \end{aligned} \quad \equiv$$

(Only if)

We prove 'only if' part by induction on the number of steps in the derivation of $A \xrightarrow[G']{*} w$

Basis If $A \xrightarrow[G']{*} w$, then $A \rightarrow w$ is in P'

By construction of P' , $A \rightarrow w$ is obtained from some production $A \rightarrow x_1 x_2 \dots x_n$ in P by erasing some (or none of the) nullable variables.

Hence $A \xrightarrow[G']{*} x_1 x_2 \dots x_n \xrightarrow[G']{*} w$.

Induction:

Assume the result for derivation in at most i steps

Let $A \xrightarrow[G']{j+1} w$

This can be split as

$A \xrightarrow[G']{*} x_1 x_2 \dots x_n \xrightarrow[G']{*} w_1 w_2 \dots w_n = w$

where $x_i \xrightarrow[G']{*} w_i$ in fewer than $j+1$ steps.

\therefore By induction hypothesis $x_i \xrightarrow[G']{*} w_i$

The first production $A \rightarrow x_1 x_2 \dots x_n$ in P' is obtained from some production $A \rightarrow \alpha$ in P by erasing some (or none of the) nullable variables in α

$$\therefore A \xrightarrow[G]{*} \alpha \xrightarrow[G]{*} x_1, x_2, \dots, x_n$$

- if $x_i \in T$ then $x_i \xrightarrow[G]{*} x_i = w_i$

if $x_i \in V$ then by induction hypothesis

$$x_i \xrightarrow[G]{*} w_i$$

$$\therefore A \xrightarrow[G]{*} x_1, x_2, \dots, x_n \xrightarrow[G]{*} w_1, w_2, \dots, w_n = w$$

Thus by the principle of induction whenever

- $A \xrightarrow[G']{*} w$, we have $A \xrightarrow[G]{*} w$ and $w \neq \epsilon$

- Applying the claim to S , we have $w \in L(G')$
iff $w \in L(G)$ and $w \neq \epsilon$
 $\Rightarrow L(G') = L(G) - \{\epsilon\}$

Corollary 1:

An algorithm to decide whether $\epsilon \in L(G)$ for a CFG:

- Construct W , the set of the nullable variables.
at most $|V|$ steps
- test whether $S \in W$

Corollary 2

If $G = (V, T, P, S)$ is a CFG, we can find an equivalent CFG $G_1 = (V_1, T, P_1, S_1)$ without null productions except $S_1 \rightarrow \epsilon$ when $\epsilon \in L(G)$. When $S_1 \rightarrow \epsilon$ is in P_1 , S_1 does not appear on the R.H.S. of any production in P_1 .

Proof:

By corollary 1, we can decide whether $\epsilon \in L(G)$

Case 1:

$$\epsilon \notin L(G)$$

Then G_1 obtained by using Theorem 4 is the required equivalent grammar.

Case 2:

$$\epsilon \in L(G)$$

Construct $G' = (V, T, P', S)$ using Theorem 4

$$\text{Then } L(G') = L(G) - \{\epsilon\}$$

$$\text{Define } G_1 = (V \cup \{S_1\}, T, P_1, S_1)$$

$$\text{where } P_1 = P' \cup \{S_1 \rightarrow S, S_1 \rightarrow \epsilon\}$$

S_1 does not appear on the R.H.S. of any production in P_1 ,

$\therefore G_1$ is the required CFG with $L(G_1) = L(G)$

Elimination of Unit productions

- A **unit production** (or a chain rule) in a CFG G is a production of the form $A \rightarrow B$ where A and B are variables in G .

Example: G is the grammar: $S \rightarrow A, A \rightarrow B, B \rightarrow C, C \rightarrow a$

$S \rightarrow A, A \rightarrow B, B \rightarrow C$, are useful just to replace S by C

Thus $L(G) = \{a\}$

If G' is the grammar $S \rightarrow a$, then $L(G') = L(G)$

Theorem 5

If G is a CFG, we can find a CFG G_1 , which has no null productions or unit productions such that

$$L(G_1) = L(G)$$

Proof:

We can apply Corollary 2 of Theorem 4 to grammar algorithm 5 construction 5 to get a CFG $G' = (V, T, P, S)$ without null productions such that

$$L(G') = L(G)$$

Let $A \in V$

Step 1: Construction of the set of variables derivable from A

Define $W_i(A)$ recursively as follows:

$$W_0(A) = \{A\}$$

$$W_{i+1}(A) = W_i(A) \cup \{B \in V \mid C \rightarrow B \text{ is in } P \text{ with } C \in W_i(A)\}$$

By definition of $W_i(A)$, $W_i(A) \subseteq W_{i+1}(A)$

As V is finite, $W_{k+1}(A) = W_k(A)$ for some $k < |V|$

$$\therefore W_{k+j}(A) = W_k(A) \quad \forall j \geq 0$$

Let $W(A) = W_k(A)$, then $W(A)$ is the set of all variables derivable from A .

Note: $B \in W(A)$ means $A \xrightarrow[G']{*} B$

Step 2: Construction of A-productions in G_1

The A-productions in G_1 are

- either i. the non-unit A-productions in G'
- or ii. $A \rightarrow \alpha$ whenever $B \rightarrow \alpha$ is in G'

Actually (ii) covers (i) as $A \in W(A)$

Now, we define $G_1 = (V, T, P, S)$ where P_1 is constructed using Step 2 for every $A \in V$

Claim: G_1 is the required grammar.

Example:

Let G be $S \rightarrow AB, A \rightarrow \alpha, B \rightarrow C|b, C \rightarrow D, D \rightarrow E, E \rightarrow \alpha$

Eliminate unit productions and get an equivalent grammar.

Soln:

Step 1: $W_0(S) = \{S\}$

$$S \rightarrow AB$$

$$W_1(S) = W_0(S) \cup \{A, \text{ if } A_2 \rightarrow A, \text{ is in } P \text{ with } A_2 \in W_0(S)\}$$

$$= W_0(S) \cup \emptyset$$

$$= \{S\}$$

$$\Rightarrow \boxed{W(S) = \{S\}}$$

Similarly $\boxed{W(A) = \{A\}}$

$$W_0(B) = \{B\}$$

$$B \rightarrow C|b$$

$$W_1(B) = W_0(B) \cup \{C\} = \{B, C\}$$

$$C \rightarrow D$$

$$W_2(B) = W_1(B) \cup \{D\} = \{B, C, D\}$$

$$D \rightarrow E$$

$$W_3(B) = W_2(B) \cup \{E\} = \{B, C, D, E\}$$

$$\boxed{W(B) = \{B, C, D, E\}}$$

$$W_0(C) = \{C\}$$

$$C \rightarrow D$$

$$W_1(C) = \{C, D\}$$

$$D \rightarrow E$$

$$W_2(C) = \{C, D, E\}$$

$$\boxed{W(C) = \{C, D, E\}}$$

$$W_0(D) = \{D\}$$

$$D \rightarrow E$$

$$W_1(D) = \{D, E\}$$

$$\boxed{W(D) = \{D, E\}}$$

Step 2:

The productions of G' are

$S \rightarrow AB, A \rightarrow a, B \rightarrow b, E \rightarrow a, B \rightarrow a, C \rightarrow a, D \rightarrow a$

(as $\{E \rightarrow a \text{ and } a \notin V\}$ is in Q with $E \in W(B)$)

$E \in W(C)$

$E \in W(D)$)

Step 3.

To complete the proof, we have to show that

$$L(G') = L(G_1)$$

- $G' = (V, T, P, S)$

$G_1 = (V, T, P_1, S)$ where P_1 consists of

i. all non-unit productions of G'

ii. productions $A \rightarrow \alpha$ whenever $B \rightarrow \alpha$ is in P
with $B \in W(A)$ and $\alpha \notin V$

i.e. if $A \rightarrow \alpha$ is in $P_1 - P$, then it is induced by $B \rightarrow \alpha$
in P with $B \in W(A)$ and $\alpha \notin V$

Now, $B \in W(A)$ implies $A \xrightarrow[G']{*} B$

$\therefore A \xrightarrow[G']{*} B \Rightarrow \alpha$

So, if $A \xrightarrow[G_1]{*} \alpha$, then $A \xrightarrow[G']{*} \alpha$

Thus,
 $L(G_1) \subseteq L(G')$

To prove $L(G') \subseteq L(G_1)$

Suppose that $w \in L(G')$ and consider a leftmost derivation of w in G' , say

$$S \xrightarrow[G']{} \alpha_1 \xrightarrow[G']{} \alpha_2 \xrightarrow[G']{} \dots \xrightarrow[G']{} \alpha_n = w$$

Let i be the smallest index such that $\alpha_i \xrightarrow[G']{} \alpha_{i+1}$ is obtained by a unit production and j be the smallest index greater than i such that $\alpha_j \xrightarrow[G']{} \alpha_{j+1}$ is obtained by a non-unit production.

So, $S \xrightarrow[G_1]{} \alpha_i$ and $\alpha_j \xrightarrow[G']{} \alpha_{j+1}$ can be written as

$$\alpha_i = w_i A_i \beta_i \Rightarrow w_i A_{i+1} \beta_i \Rightarrow$$

$$\dots \Rightarrow w_i A_j \beta_i \Rightarrow w_i f \Rightarrow \alpha_{i+1}$$

when $A_i \in W(A_i)$ and $A_j \rightarrow f$ is a non-unit production.

Hence $A_j \rightarrow f$ is a production in P_1

$\therefore \alpha_i \xrightarrow[G_1]{*} \alpha_{j+1}$ and we have $S \xrightarrow[G_1]{*} \alpha_{j+1}$

Repeating the arguments whenever some unit production occurs in the remaining part of the derivation

$$S \Rightarrow \alpha_1 \Rightarrow \dots \Rightarrow \alpha_i \Rightarrow \alpha_{i+1} \Rightarrow \dots \Rightarrow \alpha_j \Rightarrow \alpha_{j+1} \Rightarrow \dots \Rightarrow \alpha_n = w$$

We can prove that $S \xrightarrow[G_1]{*} \alpha_n = w$

Hence $L(G') \subseteq L(G_1)$

Corollary

If G is a CFG, we can construct an equivalent grammar G' which is reduced (i.e. has no useless symbols) and has no null productions or unit productions.

Proof:

We construct G' in the following way:

Algorithm 6 { Eliminate null productions to get G_1
Construction 6 { Eliminate unit productions to get G_2
 { Construct a reduced grammar G'
 equivalent to G_2

G' is the required equivalent grammar.

Note:

We have to apply the constructions only in the order given above. Otherwise we may not get the grammar in the most simplified form.