

$\phi_A(q_0, h(w))$ by defⁿ of Homomorphism.

Context Free Grammars (CFG) - Informal Comments

- * A context free grammar (CFG) is a notation for describing languages
- * It is more powerful than finite automata or RE's but still cannot define all possible languages
- * Useful for nested structures, eg. : parentheses in programming languages

Some languages which cannot be described by finite automata

$\{ 0^n 1^n \mid n \geq 1 \}$ can be described by a CFG.

$\{ 0^n 1^n 2^n \mid n \geq 1 \}$ cannot be " " " CFG.

Basic idea is to use "variables" to stand for sets of strings (i.e., languages)
~ These variables are defined recursively in terms of one another

- Recursive rules ("production") involve only concatenation
- Alternative rules for a variable allow union

- Example:

CFG for $\{0^n 1^n \mid n \geq 1\}$

$S \rightarrow 01$ "S produces 01"

$S \rightarrow 0S1$

0 01 1

Any ~~to~~ PL can be expressed in terms of CFG

gain.

$S \rightarrow$

0 0011 1.

Proof:

Basis: 01 is in the language.

Induction: If w is in the language, then so is $0w1$

CFG Formalism

- Terminals:

= symbols of the alphabet of the language being defined

- Variables • non-terminals

= a finite set of other symbols, each of which represent a language

- Start symbol

= the variables whose language is the one being defined

- Productions

A production has the form:

variable \rightarrow string of variables & terminals

Example : Define a CFG for a language consisting of all strings of a and b, in which the number of a's and the number of b's are not equal.

$L = \{abb, bba, aab, ababa, a, b, aad\}$

X starts & ends with a starts and ends with a	$S \rightarrow a S b$	$A \rightarrow a$
	$S \rightarrow b S a$	$A \rightarrow A a$
	$S \rightarrow A$	$B \rightarrow B$
	$S \rightarrow B$	$B \rightarrow B b$

Example : Formal CFG

• Here is the formal CFG for $\{0^n 1^n \mid n \geq 1\}$

Terminals = $\{0, 1\}$ $G = (V, T, P)$

Variables = $\{S\}$

Start symbol = S

Productions = $S \rightarrow 0S1$
 $S \rightarrow 01$

Derivations - Intuition.

- We derive strings in the language of CFG by
 - starting with the start symbol and repeatedly replacing some variable A by the right side of one of its productions
 - That is the "productions for A " are those that have A on the left side of the \rightarrow

$A \rightarrow \alpha$

Derivation: — Formation

- We say $\alpha A \beta \Rightarrow \alpha \gamma \beta$ if $A \rightarrow \gamma$ is a production
- Example.

$$S \Rightarrow 0S1 \Rightarrow 00S11 \Rightarrow 000111$$

$$W = 000111$$

Iterated - Derivation

- $\xRightarrow{*}$ means "zero or more derivation steps"
- Basis:

$$\alpha \xRightarrow{*} \alpha \text{ for any string } \alpha$$

- Induction:

$$\text{if } \alpha \xRightarrow{*} \beta \text{ and } \beta \Rightarrow \gamma \text{ then } \alpha \xRightarrow{*} \gamma$$

- Example: Iterated Derivation

$$S \rightarrow 0S; \quad S \rightarrow 0S1$$

$$S \xRightarrow{*} 0S1 \xRightarrow{*} 00S11 \xRightarrow{*} 000111$$

$$\text{so } S \xRightarrow{*} S; \quad S \xRightarrow{*} 0S1; \quad S \xRightarrow{*} 00S11; \\ S \xRightarrow{*} 000111$$

* Sentential Forms:

Any string of variables and/or terminals derived from the start symbol is called a sentential form.

Left most Derivations

- Say $WA \xRightarrow{lm} W\beta$ if W is a string of terminals only and $A \rightarrow \beta$ is a production.

- Also $\alpha \xRightarrow{lm} \beta$ if α becomes β by a sequence of 0 or more \Rightarrow_{lm} steps

Example: left most derivations replacing left most variable

$$S \rightarrow SS \mid (S) \mid ()$$

Derive $W = (())()$

$$S \Rightarrow (S)S \Rightarrow ((S))S \Rightarrow (())S \Rightarrow (())()$$

* Right most Derivations:

- ✓ Say $\alpha A W \xRightarrow{rm} \alpha \beta W$ if W is a string of terminals only and $A \rightarrow \beta$ is a production.

- ✓ Also $\alpha \xRightarrow{rm} \beta$ if α becomes β by a sequence of 0 or more \Rightarrow_{rm} steps.

✓ Example: Right most Derivations

✓ Balanced - parentheses grammar:

$$S \rightarrow SS \mid (S) \mid ()$$

$$W = (())()$$

$$S \Rightarrow_{rm} S(S) \Rightarrow S() \Rightarrow_{rm} (S)()$$

$$S \Rightarrow_{rm} (())()$$

Formally α is a sentential form iff $s \Rightarrow \alpha$

* Language of a Grammar

• If G is a CFG, then $L(G)$, the language of G , is $\{w \mid s \Rightarrow w\}$

• Note w must be a terminal string, s is the start symbol

• Example:

G has productions $S \rightarrow \epsilon$ and $s \rightarrow 0s1$

$L(G) = \{0^n 1^n \mid n \geq 0\}$ note: ϵ is the legitimate right side

$\{\epsilon, 01, 0011, 000111, \dots\}$

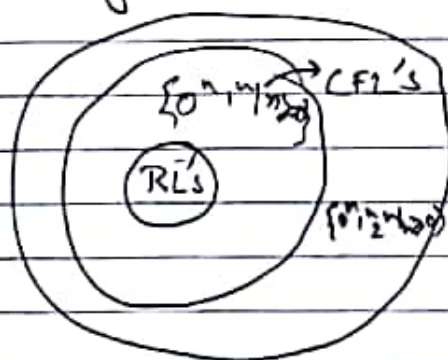
* Context - Free languages

• A language that is defined by some CFG is called a context free language.

• There are CFL's that are not regular languages, such as $\{0^n 1^n \mid n \geq 0\}$

• But not all languages are CFL's, eg. $\{0^n 1^{2n} \mid n \geq 0\}$

• Intuitively, CFL's can count two things not three



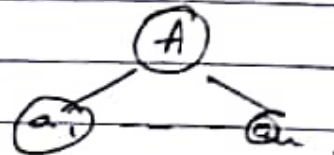
$\{0^n 1^{2n} \mid n \geq 1\}$
 $S \rightarrow S_1 1$
 $S_1 \rightarrow 0 S_1$
 $S_1 \rightarrow \epsilon$

• Proof - Part 1

✓ Induction on the height (length of the longest path from the root) of the tree

✓ Basic: height 1

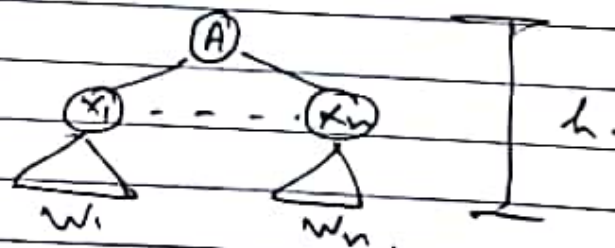
Tree looks like



$\therefore A \rightarrow a_1 \dots a_n$ must be a production.
Thus $A \Rightarrow_{\text{tree}}^* a_1 \dots a_n$

Induction:

Assume (I) for trees of height $< h$ and let this tree have height h .



By IH,

$$x_i \Rightarrow_{\text{tree}}^* w_i$$

Note: If x_i is a terminal then $x_i = w_i$

$$\text{Thus } A \Rightarrow_{\text{tree}}^* x_1 \dots x_n$$

$$\Rightarrow_{\text{tree}}^* w_1 x_2 \dots x_n$$

$$\Rightarrow_{\text{tree}}^* w_1 w_2 \dots x_n \Rightarrow_{\text{tree}}^* w_1 \dots w_n$$

$$\Rightarrow_{\text{tree}}^* w_1 \dots w_n$$

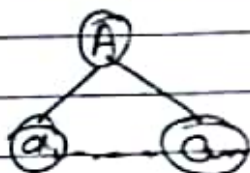
Proof : Part 2:

Given a left most derivation of a terminal string, we need to prove the existence of a parse tree.

The proof is an induction on the length of the derivation.

Basis:

If $A \Rightarrow_{\text{one step}} a_1 \dots a_n$ by a one step derivation then there must be a parse tree.



Induction:

Assume (2) for derivation of fewer than $k > 1$ steps, and let $A \Rightarrow_{\text{k steps}} w$ be a k -step derivation.

• First step is $A \Rightarrow_{\text{one step}} x_1 \dots x_n$

Now w can be divided so that first partition is derived from x_1 , the next is derived from x_2 and so on.

If x_i is a terminal then $w_i = x_i$

Induction:

That is $x_i \Rightarrow_{\text{one step}} w_i$ for all i such that x_i is a variable.

And the derivation takes fewer than k -steps.

By the IH, if x_i is a variable then there is a parse tree with root x_i and yield w_i .

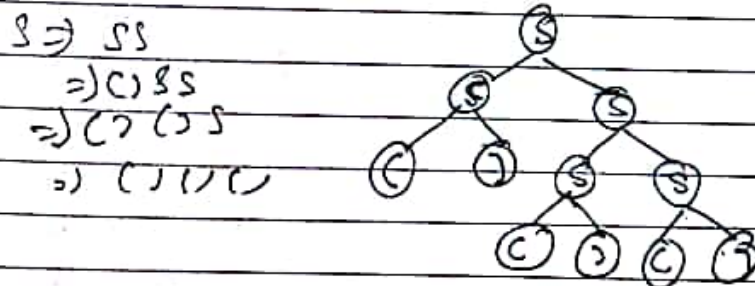
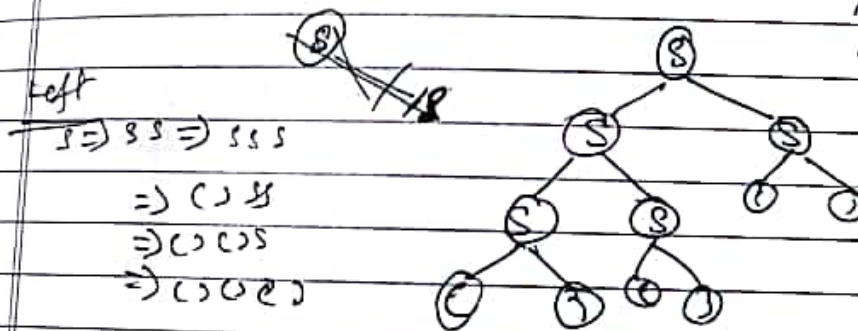
Thus there is a parse tree.

Ambiguous Grammar

* A CFG is ambiguous if there is a string in the language that is the yield of two or more parse trees.

Example: $S \rightarrow SS \mid (S) \mid ()$

Let's consider a string $()()()$ and the corresponding parse trees.



* Removing Ambiguity of a Grammar

For the balanced parentheses language, here is another CFG, which is unambiguous.

Start symbol $B \rightarrow (RB \mid \epsilon$

$R \rightarrow)RCR$

let $N = (()) ()$

$0S \Rightarrow (RB \Rightarrow ((R RB \Rightarrow (() RB$

$\Rightarrow (()) B$

$\Rightarrow (()) (RB$

$\Rightarrow (()) () B$

$\Rightarrow (()) ()$

Inherent Ambiguity

- Certain CFL's are inherently ambiguous, meaning that every grammar for the language is ambiguous.

• Example:

$\{ 0^i 1^j 2^k \mid i=j \text{ or } j=k \}$

$S \rightarrow AB \mid CD$

$A \rightarrow 0A1 \mid 01$

$B \rightarrow 2B \mid 2$

$C \rightarrow 0C \mid 0$

$D \rightarrow 1D2 \mid 12$

There are two derivations of every string with equal numbers of 0's, 1's and 2's, eg: consider 012

$S \Rightarrow AB \Rightarrow 01B \Rightarrow 012$

$S \Rightarrow CD \Rightarrow 0D \Rightarrow 012$

Variables that derive Nothing

- Consider,

$$S \rightarrow AB$$

$$A \rightarrow aA/a$$

$$B \rightarrow AB$$

Although A derives all string of a's
B derives no terminal string.

Thus S derives nothing and the language is empty

✓ Testing whether a variable derives some terminal string.

Basis: If there is a production $A \rightarrow w$, where w has no variables, then A derives a terminal string.

Induction:

If there is a production $A \rightarrow \alpha$, where α consists only of terminals and variables known to derive a terminal string, then A derives a terminal string.

Eventually we can find no more variables.

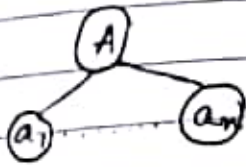
* An easy induction on the order in which variables are discovered showing that

* Conversely any variable that derives a terminal string will be discovered by this algorithm.

Proof: This proof is an induction on the height of the ~~last~~ least height parse tree by which a variable A derives a terminal string.

✓ Basis :

Height = 1 Tree looks like

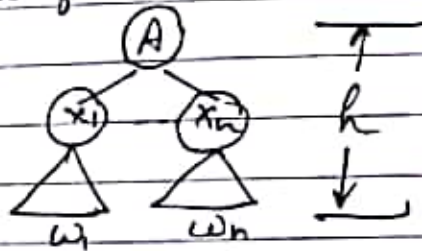


Then the Basis of the algorithm tells us that A will be discovered.

Assume Induction hypothesis for parse trees of height k and suppose A derives a terminal string via a parse tree of height h .

By IH, those x 's that are variables are discovered.

Then A will also be discovered because it has a right side of terminals and/or discovered variables.



* Algorithm to Eliminate.

① Variables that derive Nothing

Step:

1. Discover all variables that derive terminal strings
2. For all other variables, remove all productions in which they appear either on the left or the right

Example:

$S \rightarrow AB | e$

$A \rightarrow aA$

$B \rightarrow bB$

$C \rightarrow c$

Basis: A and c are identified because of $A \rightarrow c$
and $c \rightarrow c$.

Induction: S is identified because of $S \rightarrow c$.

Nothing else can be identified.

Result: $S \rightarrow c$
 $A \rightarrow aA | a$
 $c \rightarrow c$.

* Unreachable Symbol

✓ Another way a terminal or variable deserves to be eliminated is if it cannot appear in any derivation from the start symbol

✓ Basis: We can reach S (the start symbol)

✓ Induction: If we can reach A and there is a production $A \rightarrow \alpha$, then we can reach all symbols of α

When we can discover no more symbols, we have all and only the symbols that appear in derivations from S .

✓ Algorithm:

Remove from the grammar all symbols not discovered / reachable from S and all productions that involve these symbols.

Eliminating Useless Symbols

* A symbol is useful if it appears in some derivation of some terminal string from the start symbol.

* Otherwise it is useless

Eliminate all useless symbols by :

- Eliminate symbols that derive no terminal string
- Eliminate unreachable symbol

Example :

~~$S \rightarrow AB$~~

$A \rightarrow c$

$C \rightarrow c$

~~$B \rightarrow bB$~~

✓ A & C both derive terminal string

✓ B & S do not derive any terminal string.

✓ A , C and c are all unreachable symbols

Epsilon Productions (E-productions)

We can almost avoid using productions of the form $A \rightarrow \epsilon$ (called ϵ -productions)

✓ The problem is that ϵ cannot be in the language of any grammar that has no ϵ -productions

✓ Theorem :

If L is a CFL then $L - \{\epsilon\}$ has a CFG with no ϵ -productions

Nullable Symbols

• Nullable Symbols

To ~~eliminate~~ eliminate ϵ -productions, we first need to discover the nullable variables

For every nullable variable A , there is a derivation
 $A \xRightarrow{*} \epsilon$

Basis:

If there is a production $A \rightarrow \epsilon$, then A is nullable.

Induction:

If there is a production $A \rightarrow \alpha$ and all symbols of α are nullable then A is nullable.

Example:

Nullable Symbols

$S \rightarrow AB$

$A \rightarrow aA \mid \epsilon$

$B \rightarrow bB \mid \epsilon$

Basis: A is nullable because $A \rightarrow \epsilon$

Induction: B is nullable because of $B \rightarrow \epsilon$

then S is nullable because of $S \rightarrow AB$

Eliminating ϵ -productions

• Key idea:

Turn each production

$A \rightarrow x_1 \dots x_n$

into a family of productions.

- For each subset of nullable X_i 's there is one production with those eliminated from the right side "in advance"

- Except if all X_i 's are nullable do not make a production with ϵ as the right side.

* Example: Eliminating ϵ -Productions

$S \rightarrow ABC$	$B \rightarrow bB$	$S \rightarrow \epsilon$	$S \rightarrow \epsilon$
$A \rightarrow aA$	$B \rightarrow \epsilon$	$S \rightarrow AC$	$S \rightarrow \epsilon$
$A \rightarrow \epsilon$	$C \rightarrow \epsilon$	$S \rightarrow AB$	$A \rightarrow aA$
• A, B, C , and S are nullable.		$S \rightarrow C$	$B \rightarrow bB$
		$S \rightarrow A$	
		$S \rightarrow B$	

* Unit Productions:

- A unit production is one whose right side consists of exactly one variable.
- These productions can be eliminated.

Key idea:

If $A \xRightarrow{*} B$ by a series of unit production and $B \rightarrow \alpha$ is a non-unit production, then add production $A \rightarrow \alpha$.

Then drop all unit production

✓ Cleaning up a Grammar

• Theorem:

If L is a CFL then there is a CFG for $L - \{\epsilon\}$ that has:

- No useless symbols
- No ϵ -productions
- No unit production.

i.e. every right side is either a single terminal or has length ≥ 2 .

Proof: Start with a CFG for L
Perform the following steps in order:

1. Eliminate ϵ production
2. Eliminate unit productions.
3. Eliminate variables that derive no terminal string
4. Eliminate variables that derive no terminal string
5. Eliminate variables not reachable from start symbol.

* Chomsky Normal Form (CNF)

A CFG is said to be in CNF if every production is of one of these two forms

1. $A \rightarrow BC$ (right side has two variables)
2. $A \rightarrow a$ (right side has a single terminal)

Theorem:

If L is a CFL then $L - \{\epsilon\}$ has a CFG ~~in CNF~~
in CNF

Proof of CNF Theorem:

Step 1: "Clean" the grammar, so every production right side is either a single terminal or of length at least 2.

Step 2: For each right side \neq a single terminal, make the right side all variables

- For each terminal, create a new variable A_a and production $A_a \rightarrow a$
- Replace a by A_a in right side of length ≥ 2

Step 3:

Consider production $A \rightarrow BCDE$.
✓ We need variables A_c and A_e with productions
 $A_c \rightarrow c$
 $A_e \rightarrow e$
✓ Replace $A \rightarrow BCDE$ by $A \rightarrow BA_cDA_e$.

Step 4:

Break right sides longer than 2 into a ~~chain~~ chain of productions with right sides of two variables

Example: $A \rightarrow BCDE$ is replaced by

$$A \rightarrow BF$$

$$F \rightarrow CG$$

$$G \rightarrow DE$$

✓ Recall $A \rightarrow BCDE$ is replaced by
 $A \rightarrow BF$, $F \rightarrow CG$ and $G \rightarrow DE$

In the new grammar

$$A \Rightarrow BF \Rightarrow BCG \Rightarrow BCDE$$