

In general, to describe a language, there are two possible approaches:

1) recognition: describe rules (or a mechanism) to determine whether or not a certain string belongs to a language

e.g. an automaton is such a mechanism

2) generation: define rules to generate all strings of a language

A grammar is a formalism for defining a language in terms of rules that generate all strings of the language.

Since 1870, various formal methods based on the notion of rewriting or derivation have been proposed by Axel Thue, Emil Post, A.A. Markov.

In the mid 1950s the linguist Noam Chomsky introduced the notion of formal grammar with the aim of formalizing natural language. Formal grammars are in fact too simplistic to capture natural language, but they were adopted as the main formal tool to define syntactic properties of artificial languages (e.g. programming languages)

Definition: Given an alphabet Σ , a (formal) grammar G (5.2)

is a quadruple $G = (V_N, V_T, P, S)$ where

- $V_T \subseteq \Sigma$ is a finite nonempty set of symbols called terminals
- V_N is a finite nonempty set of symbols s.t. $V_N \cap \Sigma = \emptyset$, called variables or nonterminals, or syntactic categories.

Each variable represents a language

- $S \in V_N$ is called start symbol or axiom, and represents the language being defined by G

- P is a binary relation over

$$(V_N \cup V_T)^* \times V_N \times (V_N \cup V_T)^*$$

Each element $(\alpha, \beta) \in P$ is called a production or rule, and is generally written as: $\alpha \rightarrow \beta$.

Note: α ... sequence of terminals and nonterminals with at least one nonterminal

β ... sequence of terminals and nonterminals

Definition: The language $L(G)$ generated by a grammar G

is the set of strings of terminals only that can be generated starting from the axiom by a finite sequence of rule applications

Each application of a rule $\alpha \rightarrow \beta$ consists in replacing an occurrence of α with β .

Example: Palindromes:

A palindrome is a word that reads the same both forwards and backwards. (AILATIITALIA, AMOROMA)

$$L_{pal} = \{w \in \{0, 1\}^* \mid w^R = w\}$$

grammar $G_{pal} = (V_N, V_T, P, S)$, where P consists of

- 1) $S \rightarrow \epsilon$
 - 2) $S \rightarrow 0$
 - 3) $S \rightarrow 1$
- } basis: $\epsilon, 0, 1$ are palindromes

- 4) $S \rightarrow 0S0$
 - 5) $S \rightarrow 1S1$
- } induction: if S is a palindrome, $0S0$ and $1S1$ are palindromes

Example of derivation:

$$0110 : S \xrightarrow{4} 0S0 \xrightarrow{5} 01S10 \xrightarrow{1} 0110$$

$$11011 : S \xrightarrow{5} 1S1 \xrightarrow{5} 11S11 \xrightarrow{2} 11011$$

Exercise E5.1

Prove that the above grammar generates all and only palindromes over $\{0, 1\}$.

Hint: use induction on the length of the derivation

Example: natural language generation

- Sentence \rightarrow NounPhrase VerbPhrase
- NounPhrase \rightarrow Adjective NounPhrase
- NounPhrase \rightarrow Noun
- Noun \rightarrow car
- Noun \rightarrow train
- Adjective \rightarrow big
- Adjective \rightarrow broken

Notation:

1) To denote the set of productions

$$\alpha \rightarrow \beta_1, \alpha \rightarrow \beta_2, \dots, \alpha \rightarrow \beta_m$$

we use

$$\alpha \rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_m$$

2) We use $V = V_N \cup V_T$

A production of the form $\alpha \rightarrow \varepsilon$, with $\alpha \in V^* = V_N^* \cup V_T^*$ is called ε -production.

Example: $L_{eq} = \{w \in \{0,1\}^* \mid w \text{ has equal number of 0's and 1's}\}$

We have already seen that this language is not regular.

Idea to define G_{eq} s.t. $L(G_{eq}) = L_{eq}$: use induction

base: ε is in L_{eq}

induction: $\therefore 0w_A \in L_{eq}$ if w_A has one more 1 than 0

$\therefore 1w_B \in L_{eq}$ if w_B has one more 0 than 1

Characterize also languages for w_A and w_B inductively

grammar $G_{eq} = (\{S, A, B\}, \{0,1\}, P, S)$ with P

$$S \rightarrow \varepsilon \mid 0A \mid 1B$$

$$A \rightarrow 1S \mid 0AA$$

$$B \rightarrow 0S \mid 1BB$$

(A generates strings with one more 1 than 0.)

B generates strings with one more 0 than 1.)

Exercise E5.2 Prove that $L(G_{eq}) = L_{eq}$ (by induction)

Definition: Given G , the direct derivation for G is the binary relation on $(V^* \circ V_N \circ V^*) \times V^*$ defined as follows:

(φ, ψ) is in the relation if there are $\alpha \in V^* \circ V_N \circ V^*$, $\beta, \gamma, \delta \in V^*$ such that $\varphi = \gamma\alpha\delta$, $\psi = \gamma\beta\delta$ and $\alpha \rightarrow \beta \in P$.

We write $\varphi \Rightarrow \psi$ and say that ψ directly derives from φ by G .

Definition: We call derivation the reflexive, transitive closure of direct derivation. In other words, ψ derives from φ by G , written $\varphi \xRightarrow{*} \psi$ if

- a) $\varphi = \psi$, or
- b) there are $\varphi_1, \dots, \varphi_n \in V^*$ such that $\varphi_1 = \varphi$, $\varphi_n = \psi$, and $\varphi_i \Rightarrow \varphi_{i+1}$, $\forall i, 1 \leq i < n$

Definition: Given a grammar G , the language generated by G is

$$\mathcal{L}(G) = \{w \in V_T^* \mid S \xRightarrow{*} w\}$$

Notice: words in $\mathcal{L}(G)$ are constituted by terminals only.

Terminology:

- sentence: any word $w \in V_T^*$ s.t. $S \xRightarrow{*} w$, i.e. $w \in \mathcal{L}(G)$
- sentential form: any $\alpha \in V^* = (V_T \cup V_N)^*$ s.t. $S \xRightarrow{*} \alpha$

Notation: terminals: a, b, c, \dots
nonterminals: A, B, C, \dots

strings of terminals: u, v, w, x, y, z, \dots

symbols of $V = V_N \cup V_T$: X, Y, Z, \dots

sentential forms: $\alpha, \beta, \gamma, \dots$

Example: Productions for G_{eq} :

$$S \rightarrow \varepsilon \mid 0A \mid 1B$$

$$A \rightarrow 1S \mid 0AA$$

$$B \rightarrow 0S \mid 1BB$$

7/11/2008

derivation:

$$1) 001SA \Rightarrow 001S1S \quad (\text{using } A \rightarrow 1S)$$

$$2) 001S1S \Rightarrow 0011S \quad (\text{using } S \rightarrow \varepsilon)$$

$$3) 001SA \xRightarrow{*} 0011S \quad (\text{using (1) and (2)})$$

$$4) S \xRightarrow{*} 001110$$

Example: Grammar for $L_{3n} = \{a^n b^n c^n \mid n \geq 0\}$

$$G_{3n} = (\{A, B, C, S\}, \{a, b, c\}, P, S)$$

with P

$$1) S \rightarrow aSBC$$

$$2) S \rightarrow aBC$$

$$3) CB \rightarrow BC$$

$$4) aB \rightarrow ab$$

$$5) bB \rightarrow bb$$

$$6) bC \rightarrow bc$$

$$7) cC \rightarrow cc$$

} generate $a^n b^n c^n$

} moves the C's to the end

} generate the terminals from left to right

Note: we cannot simply have

$$B \rightarrow b, \quad C \rightarrow c$$

because this would generate many words not in L_{eq}

Example of derivation of $a^3 b^3 c^3$:

$$S \xRightarrow{1} aSBC \xRightarrow{2} aaSBCBC \xRightarrow{2} aaaBCBCBC$$

$$\xRightarrow{3} aaaBCBBC \xRightarrow{3} aaaBBCBBC$$

$$\xRightarrow{3} aaaBBBCCC \xRightarrow{4} aaaabBBC$$

$$\xRightarrow{5} aaaabbBBC \xRightarrow{5} aaaabbbCC$$

$$\xRightarrow{6} aaaabbbCCC \xRightarrow{7} aaaabbbccc$$

$$\xRightarrow{7} aaaabbbccc$$

Note: not each sequence of direct derivations leads to a sentence in $\mathcal{L}(G_{3n})$

e.g. with the previous grammar we could generate

$$\begin{aligned}
S &\Rightarrow \underline{a} \underline{S} BC \Rightarrow aa \underline{S} BCBC \Rightarrow aaa \underline{B} C \underline{B} C \underline{B} C \\
&\Rightarrow aaa \underline{B} C B B C C \Rightarrow aaa b \underline{C} B B C C \\
&\Rightarrow aaa b c B B C C
\end{aligned}$$

we cannot apply any other production

Also, the application of productions could go on forever (e.g. rule 1 in the previous example)

6/11/2006

Classification of Chomsky grammars into 4 groups, depending on the form of the productions:

- grammars of type 0 : no limitations
- --- 1 : context-sensitive
- --- 2 : context-free
- --- 3 : regular (or right linear)

Definition: grammar of type 0.

Productions have the most general form $\alpha \rightarrow \beta$,

with $\alpha \in V^* \circ V_N = V^*$ $\beta \in V^*$

Grammars of type 0 allow for derivations that shorten the sentential form.

A language generated by a grammar of type 0 is called of type 0.

Definition: grammar of type 1, or context sensitive

Productions have the form $\alpha \rightarrow \beta$, with

$$\alpha \in V^* \cdot V_N \cdot V^*, \quad \beta \in V^+, \quad |\alpha| \leq |\beta|$$

These productions cannot shorten the length of the sentential form to which they are applied.

A language generated by a grammar of type 1 is called of type 1, or context sensitive.

Example: G_{3n} is context sensitive. Obviously, it is also of type 0.

Definition: grammar of type 2, or context-free

Productions have the form $A \rightarrow \beta$, with $A \in V_N, \beta \in V^+$.

These productions are productions of type 1, with the additional requirement that on the left there is a single nonterminal.

A language generated by a grammar of type 2 is called of type 2, or context free.

Example $L_{2n} = \{a^n b^n \mid n \geq 1\}$ is of type 1, since the

following grammar G'_{2n} generates L_{2n}

$$S \rightarrow aB \mid SAB$$

$$BA \rightarrow AB$$

$$aA \rightarrow aa$$

$$aB \rightarrow ab$$

$$bB \rightarrow bb$$

L_{2n} is also of type 2, since it is generated by

$$S \rightarrow aSb \mid ab$$

↓ We said that grammars of type 1 are also called context-sensitive (in contrast to context-free grammar). This is justified by the original definition by Chomsky for context-sensitive grammars:

Definition: Chomsky CS-grammar

Productions have the form $\varphi_1 A \varphi_2 \rightarrow \varphi_1 B \varphi_2$
with $\varphi_1, \varphi_2 \in V^*$, $A \in V_N$, $B \in V^+$

Intuitively, A is replaced by B only if it appears "in the context" of φ_1 and φ_2 .

Theorem: Grammars of type 1 and Chomsky CS grammars generate the same class of languages

Proof: We show that, for every language L :

There is a type-1 grammar G_1 s.t. $L = \mathcal{L}(G_1)$ iff

there is a Chomsky CS grammar G_c s.t. $L = \mathcal{L}(G_c)$

"if" immediate, since each Chomsky CS grammar is of type 1

(in $\varphi_1 A \varphi_2 \rightarrow \varphi_1 B \varphi_2$ we have $B \in V^+$ and hence

$$|\varphi_1 A \varphi_2| \leq |\varphi_1 B \varphi_2|)$$

"only-if": let G_1 be a type-1 grammar for L .

We construct from G_1 a Chomsky CS grammar G_c as follows:

- 1) for each $a \in V_T$, add a new nonterminal N_a ,
- 2) replace in each production of G_1 , each $a \in V_T$ by N_a

Now all productions have the form

$$A_1 A_2 \dots A_m \rightarrow B_1 B_2 \dots B_n \text{ with } m \leq n$$

and all $A_i, B_j \in V_N$

3) For each such production $A_1 \dots A_m \rightarrow B_1 \dots B_n$, introduce a new nonterminal N , and replace the production by the following ones:

$$A_1 A_2 \dots A_m \rightarrow N A_2 \dots A_m$$

$$N A_2 \dots A_m \rightarrow N B_2 A_3 \dots A_m$$

$$N B_2 A_3 \dots A_m \rightarrow N B_2 B_3 A_4 \dots A_m$$

⋮

$$N B_2 \dots B_{m-1} A_m \rightarrow N B_2 \dots B_{m-1} B_m \dots B_n$$

$$N B_2 \dots B_m \rightarrow B_1 B_2 \dots B_n$$

(note that, due to the presence of N , these productions will not "interfere" with other ones)

Observe that all such productions are of the form

$$\gamma_1 A \gamma_2 \rightarrow \gamma_1 B \gamma_2 \quad \text{with } \gamma_1, \gamma_2 \in V^*, A \in V_N, B \in V^+$$

4) For each $a \in V_T$, add the production

$$N_a \rightarrow a \quad (\text{where } N_a \text{ is the new non-terminal associated to } a)$$

It is not difficult to see that $\mathcal{L}(G_1) = \mathcal{L}(G_c)$

(the proof is by induction on the length of the derivation of a string $w \in \mathcal{L}(G_1)$ (resp. $\mathcal{L}(G_c)$)

↑ END OPTIONAL

Definition: grammar of type 3, or regular, or right linear (5.11)

Productions have the form $A \rightarrow \delta$ with $A \in V_N$
 $\delta \in V_T \cup (V_T \circ V_N)$

(i.e., $A \rightarrow eB$ or $A \rightarrow e$, with $A, B \in V_N, e \in V_T$)

A language generated by a grammar of type 3 is called of type 3 or regular

Example: $\{e^m b \mid m \geq 0\}$ is of type 3, since it is generated by the grammar

$$\begin{aligned} S &\rightarrow eS \\ S &\rightarrow b \end{aligned}$$

Note: a grammar of type 3 is called linear, because on the right hand side of a production there is at most one non-terminal. It is called right-linear because the non-terminal is on the right of the terminal.

↓ OPTIONAL

Exercise: E 5.3: Show that grammars of type 3 generate the class of regular languages that do not contain ϵ .

Hint: given $G = (V_N, V_T, P, S)$, construct an NFA

$A_G = (V_N \cup \{F\}, V_T, \delta, S, \{F\})$ with

$B \in \delta(A, e)$ iff $A \rightarrow eB$ and

$F \in \delta(A, e)$ iff $A \rightarrow e$

Show by induction on $|w|$ that $w \in \mathcal{L}(A_G)$ iff $w \in \mathcal{L}(G)$.

Conversely, given an NFA A , construct a grammar G_A by having again nonterminals correspond to states of A .

↑ END OPTIONAL

Note on ϵ -productions (for grammars of type 1, 2, 3)

As we have defined them, grammars of type 1 (resp. 2, 3) cannot generate the empty string ϵ .

We could extend the definition by allowing also the generation of ϵ :

- if the start symbol S does not appear on the right-hand side of productions, we allow also for a production

$$S \rightarrow \epsilon \quad (\epsilon\text{-production})$$

- if the start symbol S appears on the right-hand side of productions, we introduce a new non-terminal S_{new} , make it the new start symbol, add a production

$$S_{new} \rightarrow S, \text{ and allow for } S_{new} \rightarrow \epsilon.$$

Hence, an ϵ -production used just to generate ϵ is harmless.

Note that, allowing for ϵ -productions for every non-terminal is not that harmless.

OPTIONAL
↓

Exercise: E5.4: Show that, for every language L of type 0

there is a grammar of type 1 extended with ϵ -productions on arbitrary non-terminals that generates L .

Hint: introduce a new nonterminal N_ϵ that is eliminated through an ϵ -production $N_\epsilon \rightarrow \epsilon$, and use N_ϵ to make the right-hand side of productions as long as the left-hand side.

↑ END OPTIONAL

Context-free grammars (CFGs)

11/11/2005
8/19/2006

5.13

In a CFG, the productions have the form $A \rightarrow B$ with $A \in V_N$, $B \in V^*$ (note: we allow for ϵ -productions)

Example: CFG for arithmetic expressions over variable i

$G = (\{E, T, F\}, \{i, +, *, (,), \}, P, E)$, where P is

$E \rightarrow T \mid E + T$ $E \dots$ expression

$T \rightarrow F \mid T * F$ $T \dots$ term

$F \rightarrow i \mid (E)$ $F \dots$ factor

This grammar generates, e.g., $i + i * i$

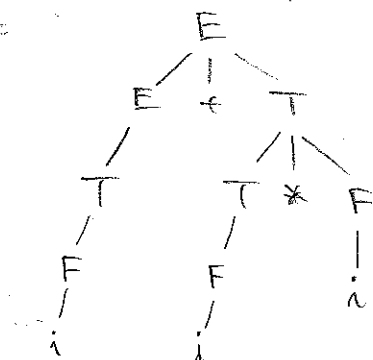
$\underline{E} \rightarrow \underline{E} + T \rightarrow \underline{T} + T \rightarrow \underline{F} + T \rightarrow i + \underline{T} \rightarrow$
 $\rightarrow i + \underline{T} * F \rightarrow i + i * \underline{F} \rightarrow i + i * i$

We can also represent a derivation of a string by a CFG by means of a tree, called parse-tree:

Is a tree whose nodes are labeled by elements of $V \cup \{\epsilon\}$ satisfying:

- 1) each interior node is labeled by a non-terminal
- 2) each leaf is labeled by a non-terminal, a terminal, or ϵ . If it is labeled by ϵ , then it is the only child of its parent
- 3) If an interior node is labeled A , and its children from left to right are labeled X_1, X_2, \dots, X_k , then there is a production $A \rightarrow X_1 X_2 \dots X_k$ in P .

Example: parse tree for $i + i * i$

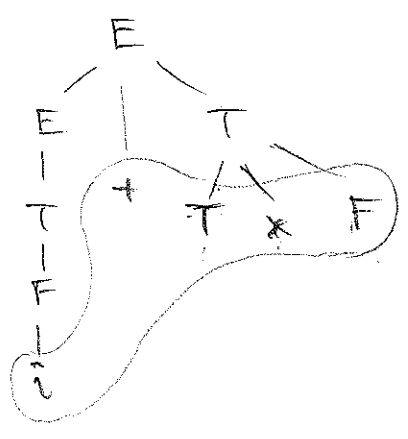


We call A-tree a subtree of the parse-tree rooted at non-terminal A.

Yield (or frontier) of a tree.

is the sequence of labels of the leaves from left to right.

Example:



Theorem: $\alpha \in V^+$ is the yield of an A-tree $\iff A \Rightarrow^* \alpha$

Proof: by induction on the height of the tree (see textbook)

Note: a parse tree does not specify a unique way to derive α from A. (the order in which non-terminals are expanded is not specified).

The parse tree specifies, however, which rule is applied for each non-terminal.

Specific derivation orders:

- leftmost derivation: obtained by traversing the tree depth-first, by first going to the left subtree, and then to the right one.

e.g. $E \xRightarrow{lm} E + T \xRightarrow{lm} i + T \xRightarrow{lm} E + T \xRightarrow{lm} i + \underline{T} \Rightarrow \dots$

- rightmost derivation: defined similarly: $E \xRightarrow{rm} E + T \Rightarrow E + \underline{T} * E$

Theorem: the following are all equivalent statements for (5.15)

a CFG $G = (V, T, P, S)$ and a string $w \in T^*$

1) $w \in \mathcal{L}(G)$ (or $S \xRightarrow{*} w$)

2) $S \xrightarrow{lm}^* w$

3) $S \xrightarrow{rm}^* w$

4) There exists an S -tree with yield w .

Proof: the equivalence of (1) and (4) follows from the previous theorem. The other equivalences are obvious.

Thus, we could always use lm -derivation as a canonical way to derive any $w \in \mathcal{L}(G)$, i.e. as a canonical way to interpret a parse tree for w .

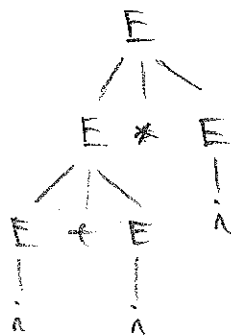
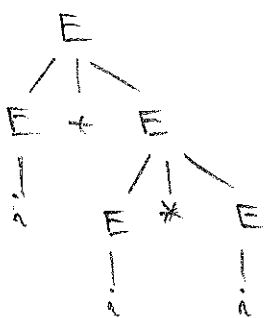
Ambiguous grammars:

$\exists w \in \mathcal{L}(G)$ could have two distinct parse trees, and hence two distinct lm -derivations

Example: another grammar for arithmetic expressions

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

$$w = i + i * i$$



These parse trees correspond to two different lm -derivations, and also to two ways of interpreting w .

Definition: A CFG G is ambiguous if for some $w \in L(G)$ there exist two distinct parse trees.

Ambiguity has to be avoided in compilers, since it corresponds to different ways of interpreting strings.

Sometimes grammar can be redesigned to remove ambiguity. (e.g., for arithmetic expressions)

This is not always possible:

Definition: A CF language is (inherently) ambiguous if all its grammars are ambiguous.

Example: $L = \{e^n b^m c^m d^m \mid n, m \geq 1\} \cup \{e^m b^m c^m d^n \mid n, m \geq 1\}$

L is CF (show for exercise)

Consider strings of the form $e^k b^k c^k d^k$.

We cannot tell whether they come from first or second types of strings in L , and any CFG must allow for both possibilities.

Properties of context-free languages

5.17

9/4/9/2008

We will study

- 1) Normal forms for CFGs (useful for proving properties of CFLs)
- 2) Expressive power \Rightarrow pumping lemma for CFLs
- 3) Closure and decision properties

Normal forms for CFGs

We look at how to simplify CFGs, while preserving the generated language.

- gain efficiency in parsing
- simplify proving properties

1) Eliminate useless symbols:

We say that $X \in V$ is useful if

$$S \Rightarrow^* \alpha X \beta \Rightarrow^* w \quad \text{with } w \in V_T^* \\ \alpha, \beta \in V^*$$

Thus, a symbol is useless (not useful) if it does not participate in any derivation of strings of the language.

\Rightarrow can be eliminated

Definition: $X \in V$ is generating if $X \Rightarrow^* w$, for $w \in V_T^*$

$X \in V$ is reachable if $S \Rightarrow^* \alpha X \beta$, for $\alpha, \beta \in V^*$

Hence, X is useful, if it is both generating and reachable

We identify useless symbols by

- 1) eliminating non-generating symbols and all their productions
- 2) --- unreachable ---

Note: it is important to do these two steps in the above order

Example: $\begin{cases} S \rightarrow AB \mid b \\ A \rightarrow \epsilon \end{cases}$ Let us consider what happens if we do first (2) and then (1)

- we eliminate unreachable symbols: all are reachable
- --- non-generating ---

we eliminate B and $S \rightarrow AB$

\Rightarrow we obtain: $\begin{cases} S \rightarrow b \\ A \rightarrow \epsilon \end{cases}$

But, if we do it in right order:

1) Eliminate non-generating symbols: B and $S \rightarrow AB$

2) --- unreachable --- : A and $A \rightarrow \epsilon$

\Rightarrow We obtain: $S \rightarrow b$

1) Eliminating non-generating symbols:

10/11/2006

Recursive algorithm to construct the set of generating symbols:

basis: mark all terminals as generating

recursive step: for each production $A \rightarrow X_1 \dots X_k$
if all of X_1, \dots, X_k are marked as generating
then mark A as generating

terminate: when no new generating symbol is found

Example: G_1 :

$$\begin{cases} S \rightarrow AB \mid AC \mid CD \\ A \rightarrow BB \\ B \rightarrow AC \mid eb \\ C \rightarrow Ce \mid CC \\ D \rightarrow Bc \mid b \mid d \end{cases}$$

$\{e, b, d\}$

$\{ \quad \rightarrow \quad, B, D \}$

$\{ \quad \quad \quad, A \}$

$\{ \quad \quad \quad, S \} \Rightarrow C$ is non-generating

\Rightarrow Remove C and all productions involving C

2) Eliminating unreachable symbols

Recursive algorithm to construct the set of reachable symbols:

basis: mark S as reachable

recursive step: for each production $A \rightarrow X_1 \dots X_n$

if A is marked as reachable

then mark X_1, \dots, X_n as reachable

terminate when no new reachable symbol is found

Example: G_2 :

$$\begin{cases} S \rightarrow AB \\ A \rightarrow BB \\ B \rightarrow eb \\ D \rightarrow b \mid d \end{cases}$$

$\{S\}$

$\{S, A, B\}$

$\{S, A, B, e, b\} \Rightarrow D, d$ are unreachable

\Rightarrow Remove D, d and all productions involving them

2) Eliminate ϵ -productions

(5.20)

ϵ -production: $A \rightarrow \epsilon$ slows down parsing

Definition: $X \in V_N$ is nullable if $X \Rightarrow^* \epsilon$

We first need to find all nullable symbols:

Recursive algorithm to construct the set of nullable symbols:

basis: if P contains $A \rightarrow \epsilon$, then mark A as nullable

inductive step: for each production $A \rightarrow X_1 \dots X_n$

if all of X_1, \dots, X_n are marked as nullable then mark A as nullable

terminate when no new nullable symbol is found

Example: G_1 :
$$\begin{cases} S \rightarrow ABC \mid BCB \\ A \rightarrow \epsilon B \mid \epsilon \\ B \rightarrow CC \mid b \\ C \rightarrow S \mid \epsilon \end{cases}$$

$\{C\}$

$\{C, B\}$

$\{C, B, S\}$

Knowing the nullable symbols allows us to compensate for the elimination of ϵ -transitions.

Example: in G_1 , since B and C are nullable, we can derive

$S \Rightarrow^* BCB$, $S \Rightarrow^* CB$, $S \Rightarrow^* BC$, $S \Rightarrow^* BB$,
 $S \Rightarrow^* C$, $S \Rightarrow^* B$, $S \Rightarrow^* \epsilon$

Hence, if we eliminate $C \rightarrow \epsilon$, we have to add direct productions for the above derivations.

Algorithm to eliminate ϵ -productions

(5.21)

1) Identify all nullable symbols

2) Replace each production $A \rightarrow X_1 \dots X_k$

by the set of all productions of the form $A \rightarrow \alpha_1 \dots \alpha_k$

where $\alpha_i = X_i$, if X_i is not nullable

$\alpha_i = X_i$ or ϵ , if X_i is nullable

3) Remove all ϵ -productions

Example: for G_1

$$\left\{ \begin{array}{l} S \rightarrow ABC \mid AB \mid AC \mid A \mid \\ \quad BCB \mid BC \mid BB \mid CB \mid B \mid C \mid \epsilon \\ A \rightarrow aB \mid \epsilon \\ B \rightarrow CC \mid C \mid \epsilon \mid b \\ C \rightarrow S \mid \epsilon \end{array} \right.$$

Finally, remove all ϵ -productions

Note: the grammar no longer generates ϵ . (this is unavoidable)

3) Eliminate unit productions

1/12/2004

Unit-production: $A \rightarrow B$ slows down parsing

Algorithm to eliminate unit-productions

1) Remove ϵ -productions

2) For all $A, B \in V_N$

if $A \Rightarrow^* B$ and $B \rightarrow \alpha$ is not unit

then add $A \rightarrow \alpha$

3) Eliminate all unit-productions

How do we find $A \Rightarrow^* B$?

Since we have no ϵ -productions: $A \Rightarrow^* B$ only if

$$A \Rightarrow B_1 \Rightarrow B_2 \Rightarrow \dots \Rightarrow B_k \Rightarrow B$$

where all B_i 's are distinct

Hence, $k \leq |V_N|$, and we can use reachability in directed graphs.

Example: $G_1 : \begin{cases} S \rightarrow A \mid B \\ A \rightarrow Se \mid e \\ B \rightarrow S \mid b \end{cases}$

reachability: $S \Rightarrow^* A, S \Rightarrow^* B$
 $B \Rightarrow^* S, B \Rightarrow^* A$

We get: $\begin{cases} S \rightarrow Se \mid e \mid b \mid A \mid B \\ A \rightarrow Se \mid a \\ B \rightarrow Se \mid e \mid b \mid S \end{cases}$

Removing unit-productions $\begin{cases} S \rightarrow Se \mid e \mid b \\ A \rightarrow Se \mid e \\ B \rightarrow Se \mid e \mid b \end{cases}$

Note: A, B are now unreachable, and hence useless

We have seen: removal of : useless symbols, ϵ -prod, unit-prod. (5.23)

Does the order of the steps matter?

Observation:

- removing useless : does not add productions at all
(and therefore not ϵ -prod. or unit-prod)
- removing ϵ -prod : could add unit-prod
- removing unit-prod :
 - needs removing ϵ -prod first
 - could create useless symbols
 - cannot create ϵ -prod.

\Rightarrow The right order for removal is

- 1) ϵ -productions
- 2) unit-productions
- 3) useless symbols : first non-generating
then unreachable

Chomsky Normal Form

13/11/2006

Definition: A CFG G is in Chomsky Normal (CNF) if all its productions are of the form

$$\begin{array}{l} A \rightarrow a \\ A \rightarrow BC \end{array} \quad \text{with} \quad \begin{array}{l} a \in V_T \\ A, B, C \in V_N \end{array}$$

Theorem: Given a CFG G with $\epsilon \notin \mathcal{L}(G)$
there is a CNF grammar G' with $\mathcal{L}(G') = \mathcal{L}(G)$.

Proof: constructive, in 3 steps

1) Eliminate ϵ -prod. and unit-prod.

\Rightarrow All productions are of the form

$$A \rightarrow \epsilon$$

$$A \rightarrow X_1 \dots X_k \quad (k \geq 2)$$

with $A \in V_N$, $\epsilon \in V_T$, $X_1, \dots, X_k \in V$

2) Remove "mixed bodies"

for each $\epsilon \in V_T$, add a new nonterminal V_ϵ and production $V_\epsilon \rightarrow \epsilon$

in each production $A \rightarrow X_1 \dots X_k$, replace ϵ with V_ϵ

\Rightarrow All productions are of the form

$$A \rightarrow \epsilon$$

$$A \rightarrow A_1 \dots A_k \quad (k \geq 2)$$

with $\epsilon \in V_T$, $A, A_1, \dots, A_k \in V_N$

3) "Factor" long productions

for each $A \rightarrow A_1 \dots A_k$ with $k \geq 3$

- add new nonterminals B_1, \dots, B_{k-2}

- replace $A \rightarrow A_1 \dots A_k$

with $A \rightarrow A_1 B_1$

$$B_1 \rightarrow A_2 B_2$$

\vdots

$$B_{k-2} \rightarrow A_{k-1} A_k$$

The grammar we get is in CNF by construction.

It is easy to show that the language is preserved

Example: G_1 $\left\{ \begin{array}{l} S \rightarrow ABB \mid \epsilon b \\ A \rightarrow B\epsilon \mid b\epsilon \\ B \rightarrow \epsilon A b B \end{array} \right.$

Step 1: nothing to do

Step 2: $\left\{ \begin{array}{l} V_\epsilon \rightarrow \epsilon \\ V_b \rightarrow b \\ S \rightarrow ABB \mid V_\epsilon V_b \\ A \rightarrow B V_\epsilon \mid V_b V_\epsilon \\ B \rightarrow V_\epsilon A V_b B \end{array} \right.$

Step 3: $\begin{array}{l} V_\epsilon \rightarrow \epsilon \\ V_b \rightarrow b \\ \left\{ \begin{array}{l} S \rightarrow A B_1 \mid V_\epsilon V_b \\ B_1 \rightarrow B B \\ A \rightarrow B V_\epsilon \mid V_b V_\epsilon \end{array} \right. \\ \left\{ \begin{array}{l} B \rightarrow V_\epsilon C_1 \\ C_1 \rightarrow A C_2 \\ C_2 \rightarrow V_b B \end{array} \right. \end{array}$