

Chomsky grammars

(5.1)

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 1820, 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

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

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

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. (AILATIDITALIA, AMOROMA)

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

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

$$\begin{array}{ll} 1) \quad S \rightarrow \epsilon & \\ 2) \quad S \rightarrow 0 & \} \text{ basis: } \epsilon, 0, 1 \text{ are palindromes} \\ 3) \quad S \rightarrow 1 & \end{array}$$

$$\begin{array}{ll} 4) \quad S \rightarrow 0 S 0 & \} \text{ induction: if } S \text{ is a palindrome,} \\ 5) \quad S \rightarrow 1 S 1 & \text{so are } 0 S 0 \text{ and } 1 S 1 \end{array}$$

Example of derivation:

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

$$11011 : S \xrightarrow{5} 1 S 1 \xrightarrow{5} 1 1 S 1 1 \xrightarrow{2} 11011$$

Exercise E 5.1 Prove that the above grammer generates all and only palindromes over $\{0,1\}$.

Hint: use induction on the length of the derivation

Example: natural language generation

Sentence \rightarrow NonPhrase VerbPhrase

NonPhrase \rightarrow Adjective NonPhrase

NonPhrase \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

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

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

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

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

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

- $w_B \in L_{eq}$ if $w_B = \dots - 0 - 1 - 0 - \dots - 1 - \dots$

Characterise also languages for w_A and w_B inductively

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

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

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

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

B generates strings with one more 0 than 1)

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

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

Definition: Given G_1 , the direct derivation for G_1 is the binary relation on $(V^* \cdot V_N \cdot V^*) \times V^*$ defined as follows:

(φ, ψ) is in the relation if there are $\alpha \in V^* \cdot V_N \cdot V^*$,

$$\alpha \in V^* \cdot V_N \cdot V^*, \quad \beta, \gamma, \delta \in V^*$$

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_1 .

Definition: We call derivation the reflexive, transitive closure of direct derivation. In other words, ψ derives from φ by G_1 , written $\varphi \xrightarrow{*} \psi$ if

a) $\varphi = \psi$, or

b) there are $\varphi_1, \dots, \varphi_n \in V^*$ such that

$$\varphi_1 = \varphi, \quad \varphi_n = \psi, \text{ and } \varphi_i \Rightarrow \varphi_{i+1}, \forall i, 1 \leq i < n$$

Definition: Given a grammar G_1 , the language generated by G_1 is

$$L(G_1) = \{w \in V_T^* \mid S \xrightarrow{*} w\}$$

Notice: words in $L(G_1)$ are constituted by terminals only.

Terminology:

- sentence: any word $w \in V_T^*$ s.t. $S \xrightarrow{*} t$, i.e. $w \in L(G_1)$

- sentential form: any $\alpha \in V^* = (V_T \cup N_F)^*$ s.t. $S \xrightarrow{*} \alpha$

Notation: terminals: s, t, 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 L_{eq} :

$$S \rightarrow \epsilon \mid OA \mid 1B$$

$$A \rightarrow 1S \mid OAA$$

$$B \rightarrow OS \mid 1BB$$

derivation:

- 1) $001SA \Rightarrow 001S1S$ (using $A \rightarrow 1S$)
- 2) $001S1S \Rightarrow 0011S$ (using $S \rightarrow \epsilon$)
- 3) $001SA \xrightarrow{*} 0011S$ (using (1) and (2))
- 4) $S \xrightarrow{*} 001110$

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

18/12/2008

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

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 $aa-aBCBC-\dots-BC$
- } moves the C's to the end
- } generate the terminals from left to right
Note: we cannot simply have $B \rightarrow b$, $C \rightarrow c$
because this would generate many words not in L_{eq}

Example of derivation of $aaabbbccc$:

$$\begin{aligned}
 S &\xrightarrow{1} a\underline{SBC} \xrightarrow{2} a\underline{aB} \underline{BCBC} \xrightarrow{2} a\underline{aa} \underline{BCBC} \underline{BC} \\
 &\xrightarrow{3} a\underline{aa} \underline{BC} \underline{BBC} \xrightarrow{3} a\underline{aa} \underline{BB} \underline{C} \underline{CC} \\
 &\xrightarrow{4} a\underline{aa} \underline{BB} \underline{B} \underline{CCC} \xrightarrow{4} a\underline{aa} \underline{bb} \underline{BB} \underline{CCC} \\
 &\xrightarrow{5} a\underline{aa} \underline{bb} \underline{B} \underline{CCC} \xrightarrow{5} a\underline{aa} \underline{bb} \underline{b} \underline{CCC} \\
 &\xrightarrow{6} a\underline{aa} \underline{bb} \underline{b} \underline{CC} \xrightarrow{6} a\underline{aa} \underline{bb} \underline{b} \underline{cc} \\
 &\xrightarrow{7} a\underline{aa} \underline{bb} \underline{b} \underline{cc}
 \end{aligned}$$

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

e.g. with the previous grammar we could generate

$$\begin{aligned} S &\Rightarrow e \underline{SBC} \Rightarrow ee \underline{SBCBC} \Rightarrow eee \underline{BC} \underline{BC} \underline{BC} \\ &\Rightarrow eee \underline{BC} BBCC \Rightarrow eee b \underline{C} BBCC \\ &\Rightarrow eee b.c BBCC \end{aligned}$$

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 $\alpha \rightarrow \beta$,

$$\text{with } \alpha \in V^0 \cup V_1 \cup V^* \quad \beta \in V^*$$

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.

The 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.

The 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}

$$S \rightarrow aB \quad | \quad SAB$$

$$BA \rightarrow AB$$

$$aA \rightarrow aa$$

$$aB \rightarrow ab$$

$$bB \rightarrow bb$$

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

$$S \rightarrow aSb \quad | \quad sb$$

OPTIONAL

↓ 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. $L = L(G_1)$ iff

there is a Chomsky CS grammar G_C s.t. $L = L(G_C)$

"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 $A_1 \cdots A_m \rightarrow B_1 \cdots B_n$, introduce a new non-terminal N , and replace the production by the following ones:

$$A_1 A_2 \cdots A_m \rightarrow N A_2 \cdots A_m$$

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

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

:

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

$$N B_2 \cdots B_n \rightarrow B_1 B_2 \cdots B_n$$

(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

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

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

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

It is not difficult to see that $\mathcal{L}(G_1) = \mathcal{L}(G_2)$

(the proof is by induction on the length of the derivation of a string $w \in \mathcal{L}(G_1)$: (resp., $\mathcal{L}(G_2)$))

\uparrow END OPTIONAL

Definition: grammar of type 3, or regular, or right linear

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 $A, B \in V_N, e \in V_T$)

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

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

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 E .

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

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

$B \in \delta(A, e) \text{ iff } A \rightarrow eB \text{ end}$

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

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

Conversely, given an NFA A , construct a grammar G_A by having each non-terminal correspond to states of A .

↑ END OPTIONAL

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

5.12

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 $S_{\text{new}} \rightarrow S$,

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

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)

In the 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

$$E \rightarrow T \mid E + T \quad E \dots \text{expression}$$

$$T \rightarrow F \mid T * F \quad T \dots \text{term}$$

$$F \rightarrow i \mid (E) \quad F \dots \text{factor}$$

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

$$E \rightarrow E + T \rightarrow T + T \rightarrow E + T \rightarrow i + T \rightarrow$$

$$\rightarrow i + T * F \rightarrow i + i * E \rightarrow i + i * i$$

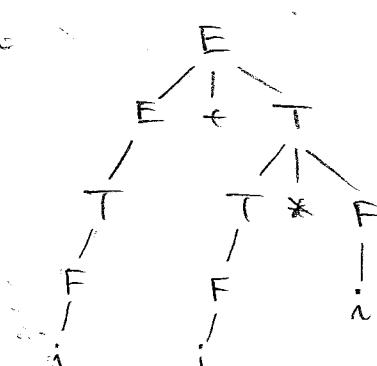
We can also represent a derivation of a string by a CFG
by means of a tree, called parse-tree:

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

$i + i * i$

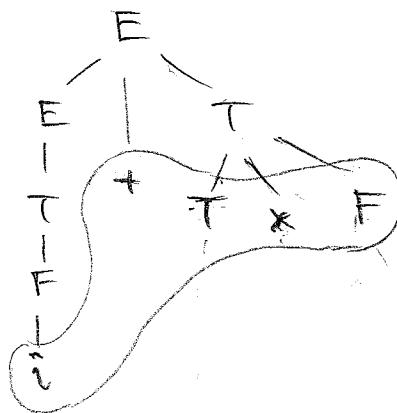


We call A-tree a subtree of the parse-tree rooted at non-terminal A.

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 $\Leftrightarrow A \Rightarrow^* \alpha$

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

8/1/2003

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. $E \xrightarrow{\text{lm}} E + T \xrightarrow{\text{lm}} I + T \xrightarrow{\text{lm}} F + T \xrightarrow{\text{lm}} i + T \xrightarrow{\text{lm}} \dots$

- rightmost derivation: defined similarly.: $E \xrightarrow{\text{rm}} E + T \xrightarrow{\text{rm}} E + T * F \xrightarrow{\text{rm}} \dots$

Theorem: the following are all equivalent statements for 5.15

- CFG $G = (V, T, P, S)$ and a string $w \in T^*$
- 1) $w \in L(G)$ (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 $w \in L(G)$; i.e. as a canonical way to interpret a parse tree for w .

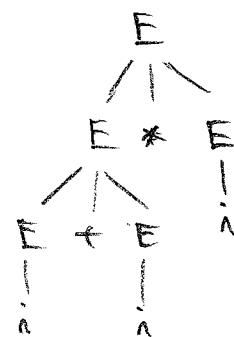
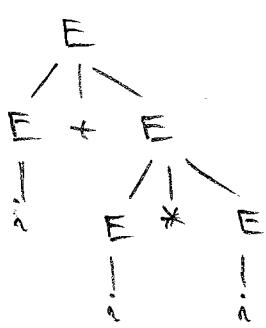
Ambiguous grammars:

• $w \in L(G)$ 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 a string.

Sometimes grammars 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

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

L is CF (show for exercise)

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

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

We will study

14/1/2008

- 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 proving
- simplify proving properties

1) Eliminate useless symbols:

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

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

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 $X \xrightarrow{*} w$, for $w \in V_T^*$

$X \in V$ is reachable if $S \xrightarrow{*} \alpha X \beta$, 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: $\begin{cases} S \rightarrow AB \mid b \\ A \rightarrow e \end{cases}$ Let us consider what happens if we do first (2) and then (1)

- we eliminate unreachable symbols: all are reachable
- ... non-generating " " : we do

we eliminate B and $S \rightarrow AB$

\Rightarrow we obtain: $S \rightarrow b$
 $A \rightarrow e$

But, if we do it in right order:

- 1) Eliminate non-generating symbols: B and $S \rightarrow AB$
- 2) ... unreachable ... : A and $A \rightarrow e$

\Rightarrow We obtain: $S \rightarrow b$

- 1) Eliminating non-generating symbols:

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: $G_1 = \left\{ \begin{array}{l} S \rightarrow AB \mid AC \mid CD \\ A \rightarrow BB \\ B \rightarrow AC \mid eb \\ C \rightarrow Ce \mid CC \\ D \rightarrow Be \mid b \mid d \end{array} \right.$

$\{e, b, d\}$

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

$\{ \dots \rightarrow \dots, A \}$

$\{ \dots \rightarrow \dots, S \} \Rightarrow C \text{ is non-generating}$

\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: $G_2 = \left\{ \begin{array}{l} S \rightarrow AB \\ A \rightarrow BB \\ B \rightarrow eb \\ D \rightarrow b \mid d \end{array} \right.$

$\{S\}$

$\{S, A, B\}$

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

\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 $X \xrightarrow{*} \epsilon$

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: for each production $A \rightarrow X_1 \dots X_n$
 step: 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: $G_1: \begin{cases} S \rightarrow ABC \mid BCB \\ A \rightarrow eB \mid e \\ B \rightarrow CC \mid b \\ C \rightarrow S \mid \epsilon \end{cases}$

15/1/2008

$\{C\}$

$\{C, B\}$

$\{C, B, S\}$

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

$$\begin{aligned} S &\xrightarrow{*} BCB, \quad S \xrightarrow{*} CB, \quad S \xrightarrow{*} BC, \quad S \xrightarrow{*} BB, \\ S &\xrightarrow{*} C, \quad S \xrightarrow{*} B, \quad S \xrightarrow{*} \epsilon \end{aligned}$$

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 $A \rightarrow X_1 \dots X_n$
by the set of all productions of the form $A \rightarrow \alpha_1 \dots \alpha_n$
where $\alpha_i = X_i$, if X_i is not nullable
 $\alpha_i = X_i \text{ or } \epsilon$, if X_i is nullable
- 3) If the resulting grammar contains $S \rightarrow \epsilon$, introduce a new start symbol S' and add the production $S' \rightarrow S|\epsilon$
- 4) Remove all ϵ -productions, except possibly the one for S' .

Example: by applying steps (1) and (2) to G_2 , we get

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

Since we have $S \rightarrow \epsilon$, we add

$$S' \rightarrow S|\epsilon$$

and remove the remaining ϵ -productions.

3) Eliminate unit productions

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 determine whether $A \Rightarrow^* B$ holds?

5.22

Since we have no ϵ -productions, we have that

$A \Rightarrow^* B$ if and only if

$A \Rightarrow B_1 \Rightarrow B_2 \Rightarrow \dots \Rightarrow B_{k-1} \Rightarrow B_k \Rightarrow B$

where all B_i 's are pairwise distinct. Hence $k \leq |V_N|$.

(if we had a sequence where two B_i 's are the same, we could eliminate all steps in between, and get a new sequence where all B_i 's are pairwise distinct)

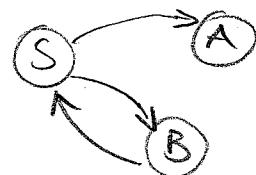
Each single derivation step $B_i \Rightarrow B_{i+1}$ must correspond to a unit production $B_i \rightarrow B_{i+1}$ of G_L .

Hence, we can detect whether $A \Rightarrow^* B$ by checking whether B is reachable from A in the graph of the unit productions:

- nodes: non-terminals

- edges: one edge $\textcircled{A} \rightarrow \textcircled{B}$ for each unit prod. $A \rightarrow B$

Example: $G_L : \begin{cases} S \rightarrow A \mid B \\ A \rightarrow Sa \mid e \\ B \rightarrow S \mid b \end{cases}$ graph of unit productions



Reachability: $S \Rightarrow^* A$ $S \Rightarrow^* B$
 $B \Rightarrow^* S$ $B \Rightarrow^* A$

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

Removing unit productions, we get

Note: A and B have become unreachable

$\begin{cases} S \rightarrow Sa \mid e \mid b \\ A \rightarrow Sa \mid e \\ B \rightarrow Sa \mid e \mid b \end{cases}$

We have seen: removal of: useless symbols, E-prod, unit-prod. 5.23

Does the order of the steps matter?

Observation:

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

⇒ The right order for removal is

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

Chomsky Normal Form

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

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

Given a CFG G , we can always construct a CFG G_C that is in CNF and such that $L(G_C) = L(G) \setminus \{\epsilon\}$.

Note: since a CFG in CNF cannot generate ϵ , if G generates ϵ , then we cannot have that $L(G_C) = L(G)$. However, apart from ϵ , the two languages are equal.

Starting from G_1 , we construct G_C in several steps:

(5.24)

- 1) Eliminate ϵ -productions (without introducing the new start symbol S' with $S' \rightarrow S|\epsilon$)
- 2) Eliminate unit-productions

\Rightarrow All productions are of the form

$$A \rightarrow a$$

$$A \rightarrow X_1 \dots X_k \quad \text{with } k \geq 2$$

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

- 3) Remove non-generating symbols

- 4) Remove unreachable symbols

- 5) Remove "mixed bodies"

for each $e \in V_T$, add a new nonterminal V_e and production $V_e \rightarrow a$

in each production $A \rightarrow X_1 \dots X_k$, replace e with V_e

\Rightarrow All productions are of the form

$$A \rightarrow e$$

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

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

- 6) "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 $A \rightarrow A_1 B_1$

$$B_1 \rightarrow A_2 B_2$$

:

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

The grammar we get is in CNF by construction.

It is easy to show that the language is preserved, except possibly for the empty string ϵ , which cannot be generated by a grammar in CNF.

Example: G₁ { $S \rightarrow ABB \mid \epsilon b$
 $A \rightarrow Be \mid b e$
 $B \rightarrow \epsilon A \& B$

Steps 4 & 5 nothing to do

Step 5: { $V_e \rightarrow \epsilon$
 $V_b \rightarrow b$
 $S \rightarrow ABB \mid V_e V_b$
 $A \rightarrow BV_e \mid V_e V_e$
 $B \rightarrow V_b A V_b B$

Step 6: $V_e \rightarrow \epsilon$
 $V_b \rightarrow b$
{ $S \rightarrow AB_1 \mid V_e V_b$
 $B_n \rightarrow BB$
 $A \rightarrow BV_e \mid V_b V_e$
{ $B \rightarrow V_e C_1$
 $C_1 \rightarrow AC_2$
 $C_2 \rightarrow V_b B$

Note: If the original grammar generated ϵ , and we went that the grammar in CNF also generates ϵ , we can execute step 1 by introducing the new start symbol S' and the productions: $S' \rightarrow S \mid \epsilon$

Step 2 will then replace the unit production $S' \rightarrow S$ by other productions, but none of the transformations 2-5 will introduce a production where S' is on the right side.

Therefore, in the end, we will have a grammar in CNF, apart for the production $S' \rightarrow \epsilon$.