

Context free grammar $G = (N, \Sigma, P, S)$

Parse tree is a tree satisfying the following conditions.

1. Each interior node is labelled with an element of N
2. Each leaf node is labelled with Σ or ϵ .
3. if an interior node is labelled A and its children are labelled B_1, \dots, B_k .

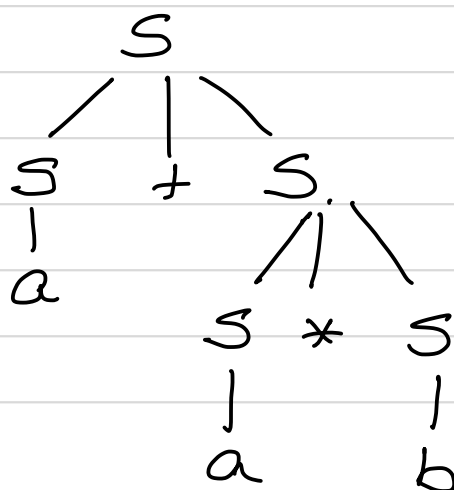
Then $A \rightarrow B_1 B_2, \dots, B_k \in P$.

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

$$I \rightarrow a \mid b$$

String $a + a * b$

$$S \rightarrow S + S \rightarrow S + S * S \xrightarrow{*} a + a * b.$$



$$S \xrightarrow{G_1} SS \xrightarrow{G_1} OS \xrightarrow{G_1} 01$$

$$G = (N, \Sigma, P, S)$$

CNF

$$S \rightarrow SS \mid 0 \mid 1 \mid \epsilon.$$

$$S \rightarrow SS \mid 0 \mid 1$$

$$S \xrightarrow{G_1} SS \xrightarrow{G_1} SSS \xrightarrow{G_1} SSSS \xrightarrow{G_1} SSSS \xrightarrow{G_1} SSS \xrightarrow{G_1} SS \rightarrow OS \rightarrow 01$$

Chomsky Normal Form. (CNF)

$G = (N, \Sigma, P, S)$ is in CNF if all productions

$$A \rightarrow BC \quad A \rightarrow a \quad A, B, C \in N, a \in \Sigma.$$

Example. $S \rightarrow [S] \mid SS \mid \epsilon \Rightarrow G_1$

$G_2 \downarrow$

$$S \rightarrow AB \mid AC \mid SS \quad C \rightarrow SB \quad A \rightarrow [\quad B \rightarrow]$$

$$L(G_1) = L(G_2).$$

One step progress. $\begin{cases} - \# \text{ of Nonterminals increase by 1} \\ - \# \text{ of terminal increase by 1.} \end{cases}$

Theorem. For any CFG G , there is a CFG G' in Chomsky normal form st $L(G') = L(G) - \{\epsilon\}$.

Lemma 1. For any CFG $G = (N, \Sigma, P, S)$ there is a CFG G' with no ϵ -productions or unit-productions such that $L(G') = L(G) - \{\epsilon\}$

Proof. Let \hat{P} be the smallest set of productions containing P and closed under the rules:

(a) if $A \rightarrow \alpha B \beta$ and $B \rightarrow \epsilon$ are in \hat{P} then $A \rightarrow \alpha \beta \in \hat{P}$

(b) if $A \rightarrow B$ and $B \rightarrow \gamma$ are in \hat{P} then $A \rightarrow \gamma \in \hat{P}$.

Note: \hat{P} is finite \rightarrow

$\left\{ \begin{array}{l} \text{Finitely many new production rules are added} \\ \text{Each new RHS is a substring of an old RHS.} \end{array} \right\}$

$$\hat{G} = (N, \Sigma, \hat{P}, S)$$

We have $L(G) \subseteq L(\hat{G})$ Since $P \subseteq \hat{P}$

$L(G) = L(\hat{G})$ - each new production was

included because of rule (a) or (b) - can be

simulated in 2 steps by two productions that caused it to be included.

Claim 2. For any non-null $x \in \Sigma^*$, any derivation

$S \xrightarrow[\hat{G}]^* x$ of minimum length does not use ϵ -or unit productions.

Proof. Let $x \neq \epsilon$. Let $S \xrightarrow[\hat{G}]^* x$ be the minimum length derivation.

Suppose an ϵ -production $B \rightarrow \epsilon$ is used at some point

$$S \xrightarrow[\hat{G}]^* \gamma B S \xrightarrow[\hat{G}]^1 \gamma S \xrightarrow[\hat{G}]^* x.$$

At least one of γ or S is non-null $\Rightarrow B$ was introduced from a production of the form $A \rightarrow \alpha B \beta$.

$$S \xrightarrow[\hat{G}]^m \gamma A \theta \xrightarrow[\hat{G}]^1 \gamma \alpha B \beta \theta \xrightarrow[\hat{G}]^n \gamma B S \xrightarrow[\hat{G}]^1 \gamma S \xrightarrow[\hat{G}]^k x$$

for $m, n, k \geq 0$

By rule (a) $A \rightarrow \alpha \beta \in \hat{P}$.

But then we have a strictly shorter derivation of x

$$S \xrightarrow[\hat{G}]^m \gamma A \theta \xrightarrow[\hat{G}]^1 \gamma \alpha \beta \theta \xrightarrow[\hat{G}]^n \gamma S \xrightarrow[\hat{G}]^k x.$$

This gives a contradiction.

Unit Productions

Let $x \neq \epsilon$. Consider a derivation $S \xrightarrow[\hat{G}]^* x$ of minimum length.

Suppose a unit production $A \rightarrow B$ is used at some point

$$S \xrightarrow[\hat{G}]^* \alpha A \beta \xrightarrow[\hat{G}]^1 \alpha B \beta \xrightarrow[\hat{G}]^* x.$$

B must be removed later by applying a production $B \rightarrow \gamma$.

$$S \xrightarrow[\hat{G}]^m \alpha A \beta \xrightarrow[\hat{G}]^1 \alpha B \beta \xrightarrow[\hat{G}]^n \eta B \theta \xrightarrow[\hat{G}]^1 \eta \gamma \theta \xrightarrow[\hat{G}]^k x.$$

By rule (b), $A \rightarrow \gamma \in \hat{P}$.

But then, there is a shorter derivation of x

$$S \xrightarrow[\hat{G}]^m \alpha A \beta \xrightarrow[\hat{G}]^1 \alpha \gamma \beta \xrightarrow[\hat{G}]^n \eta \gamma \theta \xrightarrow[\hat{G}]^k x.$$

This is a contradiction.

Claim 2 implies we can remove the ϵ -productions and unit productions from \hat{P} without changing the language.

Chomsky Normal Form.

By Lemma 1, $L(G) = L(\hat{G})$ and \hat{P} does not have ϵ -productions or unit productions.

For each terminal $a \in \Sigma$ introduce a new nonterminal A_a and add the production rule $A_a \rightarrow a$.

Replace all occurrences of a on the RHS of old productions (except productions of the form $B \rightarrow a$) with A_a . Then all productions are of the form:

$$A \rightarrow a \quad \text{or} \quad A \rightarrow \underbrace{B_1 B_2 \dots B_k}_{\text{nonterminals}} \quad k \geq 2$$

For any production of the form $A \rightarrow B_1 B_2 \dots B_k$ with $k \geq 3$, introduce a new nonterminal C and replace with

$$A \rightarrow B_1 C \quad \text{and} \quad C \rightarrow B_2 \dots B_k.$$

Repeat until all RHS of all productions are of length at most 2.

(a) if $A \rightarrow \alpha B \beta$ and $B \rightarrow \epsilon$ are in \hat{P} then $A \rightarrow \alpha \beta \in \hat{P}$

(b) if $A \rightarrow B$ and $B \rightarrow \gamma$ are in \hat{P} then $A \rightarrow \gamma \in \hat{P}$.

Example 1: $\{a^n b^n \mid n \geq 0\} - \{\epsilon\} = \{a^n b^n \mid n \geq 1\}$.

$$S \rightarrow a S b \mid \epsilon$$

$$S \rightarrow a S b \mid a b$$

Add nonterminals A, B

$$S \rightarrow A S B \mid A B \quad A \rightarrow a \quad B \rightarrow b$$

Add nonterminal C , replace $S \rightarrow A S B$ with

$$G': S \rightarrow A B \mid A C, C \rightarrow S B, A \rightarrow a, B \rightarrow b$$

Balanced Parenthesis $S \rightarrow [S] \mid S S \mid \epsilon$

$$S \rightarrow [S] \mid S S \mid []$$

Add new nonterminals A, B

$$S \rightarrow A S B \mid S S \mid A B, \quad A \rightarrow [, B \rightarrow]$$

Add a new nonterminal C . Replace $S \rightarrow A S B$ with $S \rightarrow A C$ and $C \rightarrow S B$.

$$G': S \rightarrow A B \mid A C \mid S S, \quad C \rightarrow S B, \quad A \rightarrow [, B \rightarrow]$$