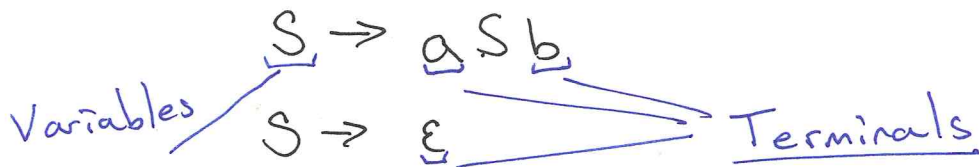


Context-Free Grammars

$$A = \{ a^n b^n \mid n \geq 0 \}$$

A grammar is a set of substitution rules, called productions



The terminals are the elements of the alphabet Σ

Grammars are used to describe a language by generating all strings of the language.

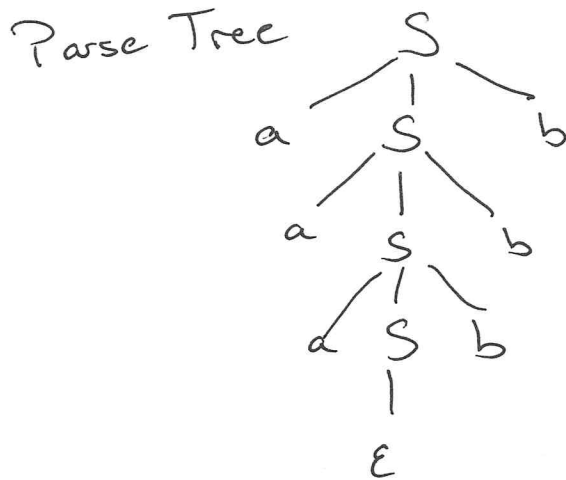
1. Write down the start variable
- left-hand side of the first rule unless otherwise specified.
2. Find a variable and a rule that starts with that variable. Replace that variable with the right hand side of the rule.
3. Repeat step 2 until no variables remain

This sequence of substitutions to obtain a string is called a derivation

$$w = aaabbb$$

$$S \Rightarrow aSb \Rightarrow aaSbb \Rightarrow aaasbbb \Rightarrow aaabbb$$

yields



All strings generated in this way constitute the language of the grammar

$$L(G)$$

$$S \rightarrow aSb \mid \epsilon$$

Formal Definition of a Context-Free Grammar

$$G = (V, \Sigma, R, S)$$

V = a finite set of variables

Σ = a finite set of terminals

R = a finite set of rules

$S \in V$ is the start variable

\Rightarrow yields (one step)

$*\Rightarrow$ derives (one or more steps)

$\{w \mid w \text{ is a palindrome}\} \quad \Sigma = \{a, b\}$

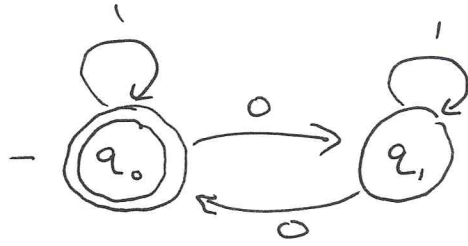
$$A \rightarrow aAa$$

$$A \rightarrow bAb$$

$$A \rightarrow a \mid b \mid \epsilon$$

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

$\{w \mid w \text{ is a set of properly nested parentheses}\}$



$\{w \mid w \text{ contains an even number of } 0\text{'s}\}$

$A \rightarrow 0B$

$A \rightarrow 1A$

$B \rightarrow 0A$

$B \rightarrow 1B$

$A \rightarrow \epsilon$

Context-free grammars that generate
Regular Languages

A context-free language is any language that can be generated with a context-free grammar.

Theorem

The class of context-free languages are closed under union.

Proof by construction:

Let A and B be context-free languages.

Let $G_1 = (V_1, \Sigma, R_1, S_1)$ s.t. $L(G_1) = A$

Let $G_2 = (V_2, \Sigma, R_2, S_2)$ s.t. $L(G_2) = B$

Construct G to ~~recognize~~^{generate} $A \cup B$

$$G = (V, \Sigma, R, S)$$

$$V = V_1 \cup V_2 \cup \{S\}$$

$$R = R_1 \cup R_2 \cup \{S \rightarrow S_1 \mid S_2\}$$

$$A = \{a^n b^n \mid n \geq 0\}$$

$$B = \{a^{2^n} \mid n \geq 0\}$$

$$G_1 \quad S_1 \rightarrow aS_1, b \mid \epsilon$$

$$G_2 \quad S_2 \rightarrow aaS_2 \mid \epsilon$$

$$G \quad S \rightarrow S_1 \mid S_2$$

Does this help us to find if
context-free languages are closed
under intersection?

No

Theorem

The class of context-free languages is closed under concatenation. For any context-free languages A, B $A \circ B$ is context-free.
Proof by construction

$$S \rightarrow AB$$

$$A \rightarrow aAb$$

$$A \rightarrow \epsilon$$

$$B \rightarrow aBb \mid bBa \mid BB \mid \epsilon$$

Let $G_1 = (V_1, \Sigma, R_1, S_1)$ and $L(G_1) = A$

Let $G_2 = (V_2, \Sigma, R_2, S_2)$ and $L(G_2) = B$

Construct $G = (V, \Sigma, R, S)$ to generate $A \circ B$

$$V = V_1 \cup V_2 \cup \{S\}$$

$$R = R_1 \cup R_2 \cup \{S \rightarrow S_1 S_2\}$$

Theorem

The class of context-free languages are closed under Kleene Star. If A is context-free then A^* is context-free.

Proof by Construction

$$\{a^n b^n \mid n \geq 0\}^*$$

$$S \rightarrow aSb \mid \epsilon$$

new start
state

$$\rightarrow T \rightarrow ST \mid \epsilon$$

Let $G = (V, \Sigma, R, S)$ where $L(G) = A$

Construct $G' = (V', \Sigma, R', S')$ to generate A^*

$$V' = V \cup \{S'\}$$

$$R' = R \cup \left\{ \begin{array}{l} S' \rightarrow SS' \\ S' \rightarrow \epsilon \end{array} \right\}$$

Ambiguity

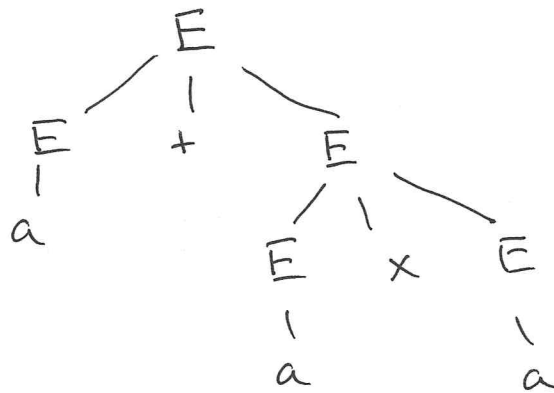
$$E \rightarrow E + E \mid E - E \mid$$

$$E \times E \mid E \div E \mid$$

$$(E) \mid a$$

$$s = a + a \times a$$

$$E \Rightarrow E + E \Rightarrow a + E \Rightarrow a + E \times E \Rightarrow a + a \times E \\ \Rightarrow a + a \times a$$



$$S = a + a \times a$$

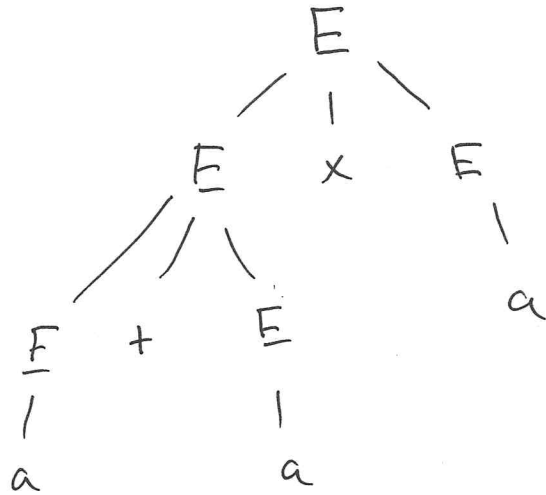
$$E \Rightarrow E \times E$$

$$\Rightarrow E + E \times E$$

$$\Rightarrow a + E \times E$$

$$\Rightarrow a + a \times E$$

$$\Rightarrow a + a \times a$$



Def: If a string s can be derived in several different ways then that string is derived ambiguously in that grammar.

If a grammar generates some string ambiguously then that grammar is ambiguous

Why is Ambiguity a Problem?

Why didn't we discuss ambiguity with regards to regular languages?

- with regular languages we don't really care about assigning internal structure to the strings.

Why are context-free languages different?

- The internal structure is often important with a string from a context-free language. And when we want to assign a meaning to a string, we almost always want to assign a unique meaning.

Changing EXPR grammar to take into account precedence rules.

$$E \rightarrow E + T \mid E - T \mid T$$

$$T \rightarrow T \times F \mid T \div F \mid F$$

$$F \rightarrow (E) \mid a$$

$a + a \times a$

$$E \Rightarrow E + T \Rightarrow T + T \Rightarrow F + T$$

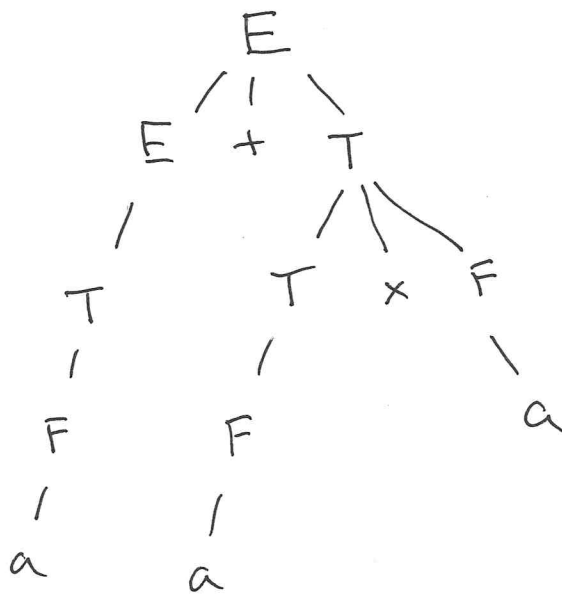
$$\Rightarrow a + T$$

$$\Rightarrow a + T \times F$$

$$\Rightarrow a + F \times F$$

$$\Rightarrow a + a \times F$$

$$\Rightarrow a + a \times a$$



Def

A context-free language that can be generated only by an ambiguous grammar is called inherently ambiguous.

$$\{a^i b^j c^k \mid i=j \text{ or } j=k\}$$

$\langle \text{expr} \rangle ::= \text{true} \mid \text{false}$

$\langle \text{stmt} \rangle ::= \text{if } \langle \text{expr} \rangle \text{ then } \langle \text{stmt_list} \rangle$
 $\mid \text{if } \langle \text{expr} \rangle \text{ then } \langle \text{stmt_list} \rangle \text{ else } \langle \text{stmt_list} \rangle$

$\langle \text{stmt_list} \rangle ::= \langle \text{stmt} \rangle$

...

if true then

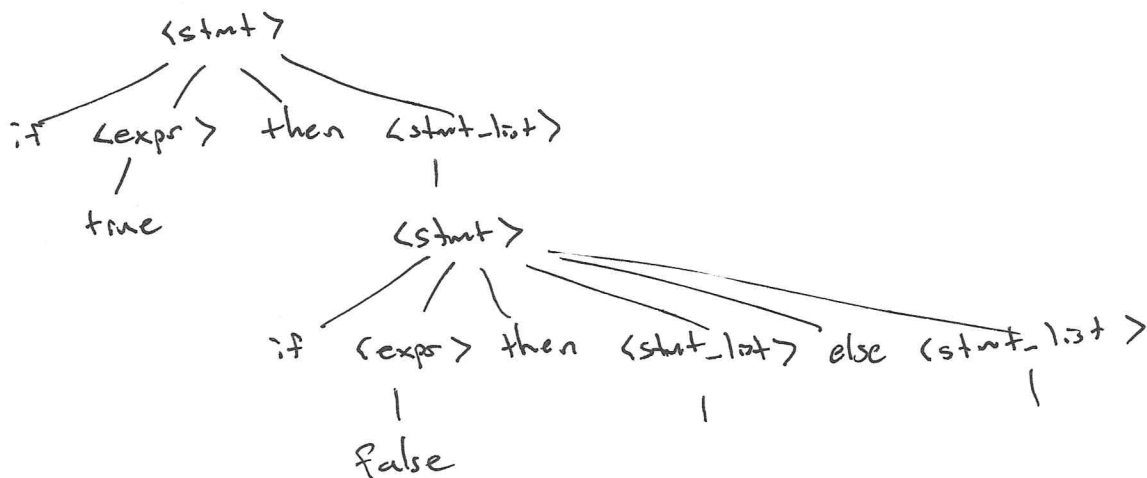
if false then

else

if true then

if false then

else



Chomsky Normal Form

In a CNF grammar $G = (V, \Sigma, R, S)$ all rules have one of the following two forms.

- $X \rightarrow a$ where $a \in \Sigma$
- $X \rightarrow BC$ where $B, C \in V$

All parse trees generated from CNF grammars have a branching factor of two.

which makes the functionally binary trees.

Why is this useful?

- Parsers can exploit efficient data structures for storing and manipulating binary trees.
- Every derivation of a string w contains $|w|-1$ applications of some rule of the form $X \rightarrow BC$ and $|w|$ applications of some rule of the form $X \rightarrow a$

This makes it much easier to define procedure to determine whether w can be generated by a CNF grammar G .

Chomsky Normal Form

In a CNF grammar all rules have one of the following two forms

$$X \rightarrow a$$

$$X \rightarrow BC \quad (\text{B and C cannot be the start variable})$$

IF X is the start variable we also permit

$$X \rightarrow \epsilon$$

Theorem

Any context-free language is generated by a context-free grammar in Chomsky normal form.

Proof idea

Convert an arbitrary grammar into Chomsky normal form.

Several stages

At each stage the grammar must generate the same language.

For some grammar G

1. Apply some transformation to G that eliminates undesirable property 1. Show that the language generated by G is unchanged.
2. Apply another transformation to G that eliminates undesirable property 2. Show that the language generated by G is unchanged and...
that undesirable property 1 has not been reintroduced.
3. Continue until the grammar has the desired form.

Undesirable properties

→ add new start variable

1. ϵ rules ($A \rightarrow \epsilon$ where $A \neq S$)

2. remove all unit rules

3. convert all remaining rules to the proper form.

- remove mixed rules

Example 0

$$S \rightarrow aSb$$

$$S \rightarrow \epsilon$$

$$\begin{aligned} \textcircled{1} \quad S_0 &\rightarrow S \\ S &\rightarrow aSb \\ S &\rightarrow \epsilon \end{aligned}$$

$$\textcircled{2} \quad S_0 \rightarrow S \mid \epsilon$$

$$S \rightarrow aSb \mid ab$$

$$\textcircled{3} \quad S_0 \rightarrow aSb \mid ab \mid \epsilon$$

$$S \rightarrow aSb \mid ab$$

④

$$S_0 \rightarrow U_a S U_b \mid U_a U_b \mid \epsilon$$

$$S \rightarrow U_a S U_b \mid U_a U_b$$

$$U_a \rightarrow a$$

$$U_b \rightarrow b$$

$$U_a T$$

$$T \rightarrow S U_b$$

⑤

$$S_0 \rightarrow U_a T \mid U_a U_b \mid \epsilon$$

$$T \rightarrow S U_b$$

$$S \rightarrow U_a T \mid U_a U_b$$

$$U_a \rightarrow a$$

$$U_b \rightarrow b$$

example 3

$$S \rightarrow RaRaR$$

$$R \rightarrow aR \mid bR \mid \epsilon$$

ϵ -rule removal

$$S \rightarrow RaRaR \mid aRaR \mid RaRa \mid RaR \mid Ra \mid aRa \mid aR \mid a$$

$$R \rightarrow aR \mid bR \mid a \mid b$$

mit rules ✓

$$S \rightarrow R \overset{\checkmark}{u_a} R \overset{\checkmark}{u_a} R \mid \overset{\checkmark}{u_a} R \overset{\checkmark}{u_a} R \mid R \overset{\checkmark}{u_a} R \overset{\checkmark}{u_a} \mid R \overset{\checkmark}{u_b} R \overset{\checkmark}{u_b} \mid R \overset{\checkmark}{u_a} \overset{\checkmark}{u_a} \mid \overset{\checkmark}{u_a} R \overset{\checkmark}{u_a} \mid \overset{\checkmark}{u_b} \overset{\checkmark}{u_a} R \mid \overset{\checkmark}{u_a} \overset{\checkmark}{u_a}$$

$$R \rightarrow \overset{\checkmark}{u_a} R \mid \overset{\checkmark}{u_b} R \mid a \mid b$$

$$u_a \rightarrow a$$

$$u_b \rightarrow b$$

remove long

$$S \rightarrow RS_1 \mid u_a S_2 \mid u_b u_b$$

$$R \rightarrow u_a R$$

$$S_1 \rightarrow u_a S_3 \mid u_a u_a$$

$$R \rightarrow u_b R \mid a \mid b$$

$$S_2 \rightarrow RS_4 \mid R u_a \mid u_a R$$

$$u_a \rightarrow a$$

$$S_3 \rightarrow RS_5 \mid R u_a \mid u_a R$$

$$u_b \rightarrow b$$

$$S_4 \rightarrow u_a R$$

$$S_5 \rightarrow u_a R$$

Example 1

$$A \rightarrow aBCa$$

$$B \rightarrow C|a$$

$$C \rightarrow D|b$$

$$D \rightarrow bD|E$$

add new start variable

$$S \rightarrow A$$

$$A \rightarrow aBCa$$

$$B \rightarrow C|a$$

$$C \rightarrow D|b$$

$$D \rightarrow bD|E$$

remove ϵ -rules

$$S \rightarrow A$$

$$A \rightarrow aBCa$$

$$B \rightarrow C|a$$

$$C \rightarrow D|b|\epsilon$$

$$D \rightarrow bD|b$$

$$S \rightarrow A$$

$$A \rightarrow aBCa | aBa$$

$$B \rightarrow C|a|\epsilon$$

$$C \rightarrow D|b$$

$$D \rightarrow bD|b$$

$$S \rightarrow A$$

$$A \rightarrow aBCa | aBa | aCa | aa$$

$$B \rightarrow C|a$$

$$C \rightarrow D|b$$

$$D \rightarrow bD|b$$

remove mit rules

$$S \rightarrow aBCa \mid aBa \mid aCa \mid aa$$

$$B \rightarrow D \mid b \mid a$$

$$C \rightarrow bD \mid b \mid b$$

$$D \rightarrow bD \mid b$$

$$S \rightarrow aBCa \mid aBa \mid aCa \mid aa$$

$$B \rightarrow bD \mid b \mid a$$

$$C \rightarrow bD \mid b$$

$$D \rightarrow bD \mid b$$

remove mixed

$$S \rightarrow TBCT \mid TBT \mid TCT \mid TT$$

$$B \rightarrow RD \mid b \mid a$$

$$C \rightarrow RD \mid b$$

$$D \rightarrow RD \mid b$$

$$T \rightarrow a$$

$$R \rightarrow b$$

Remove long

$$S \rightarrow TS_1 \mid TS_2 \mid TS_3 \mid TT$$

$$S_1 \rightarrow BS_4$$

$$S_4 \rightarrow CT$$

$$S_2 \rightarrow BT$$

$$S_3 \rightarrow CT$$

$$B \rightarrow RD \mid b \mid a$$

$$C \rightarrow RD \mid b$$

$$D \rightarrow RD \mid b$$

$$T \rightarrow a$$

$$R \rightarrow b$$