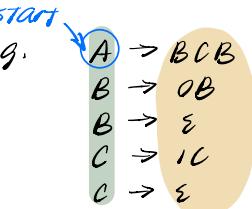


Context-free Grammar (CFG)

A CFG consists of the following components:

- A set of variables usually denoted by capital letters e.g. S, A, B, C, \dots
- An alphabet : The alphabet letters are called terminals e.g. $\Sigma = \{0, 1\}$
- A set of rules of the form e.g.



- A start variable
(usually the LHS of the first row)

use CFG to generate a language

- ① Start with the start variable
- ② A rule of the form $var \rightarrow w$ tells us that we can take one of the occurrences of the var and replace it with w
- ③ We're done when the string has no variable

e.g. $A \Rightarrow BCB \Rightarrow 0BCB \Rightarrow 00BCB \Rightarrow 00BC \Rightarrow 00B1C \Rightarrow 0010 \Rightarrow 001$

$\Rightarrow A \in V \text{ start.}$

Notation A CFG is a 4-tuple (V, Σ, R, A)

$\stackrel{b}{\Rightarrow}$
 $s \Rightarrow^b w \text{ where}$
 $s \in V \text{ and } w \in (\Sigma \cup V)^*$

- $w \Rightarrow u$ if u can be obtained from w by applying one rule.
- $w \Rightarrow^* u$ if u can be obtained from w by applying a sequence of rules

DEF The language of a CFG with a start variable A is

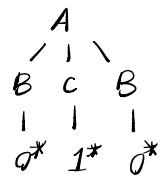
$$\{w \in \Sigma^* \mid A \Rightarrow^* w\}$$

Example

$$\begin{aligned} ① \quad & A \Rightarrow 0A1 \\ & A \Rightarrow \epsilon \end{aligned}$$

$$\Rightarrow \{0^n 1^n \mid n \geq 0\}$$

$$\begin{array}{ll}
 \textcircled{2} \quad A \Rightarrow BCB & \{ w \in \Sigma^* \mid B \Rightarrow^* w \} = 0^* \\
 B \Rightarrow 0B & \{ w \in \Sigma^* \mid C \Rightarrow^* w \} = 1^* \\
 B \Rightarrow \epsilon & \\
 C \Rightarrow 1C & \Rightarrow 0^* 1^* 0^* \\
 C \Rightarrow \epsilon &
 \end{array}$$



\textcircled{3} Design a CFG whose language is
 $\{ w \mid w \text{ starts and ends with the same symbol}, |w| \geq 3 \}$ over $\Sigma = \{0, 1\}$

$$\begin{array}{ll}
 A \Rightarrow 0B0 & B \Rightarrow \epsilon \\
 A \Rightarrow 1B1 & B \Rightarrow 0B \\
 & B \Rightarrow 1B
 \end{array}$$

\textcircled{4} Palindromes = $\{ w \mid w = w^R \}$ over $\Sigma = \{0, 1\}$

$$\begin{array}{ll}
 A \Rightarrow 0AO & A \Rightarrow 0 \\
 A \Rightarrow 1AI & A \Rightarrow 1 \\
 & A \Rightarrow \epsilon
 \end{array}$$

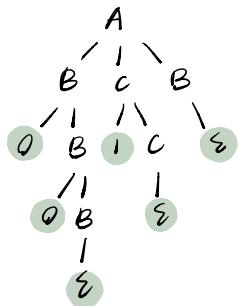
\textcircled{5} $\{ w \mid \text{number of } 1 \text{ is divisible by } 2 \}$

$A \Rightarrow AIAIA \dots IA$ The first rule can be applied repeatedly to generate an even number of A's

Then we can use $A \Rightarrow^* 0^*$ to replace those A with 0's.

Parsing tree

$$A \Rightarrow BCB \Rightarrow 0BCB \Rightarrow 00BCB \Rightarrow 00BC \Rightarrow 00B1C \Rightarrow 0010 \Rightarrow 001$$



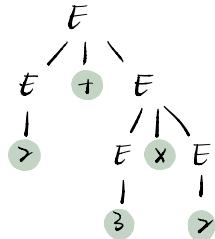
Notation
$$\left. \begin{array}{l} A \Rightarrow BC \\ A \Rightarrow \epsilon \\ A \Rightarrow AB1 \end{array} \right\} \Leftrightarrow A \Rightarrow BC \mid \epsilon \mid AB1$$

Consider the following CFA $\Sigma = \{+, \times, 0, \dots, 9\}$

$$E \Rightarrow E+E \mid E \times E \mid 0 \mid 1 \mid \dots \mid 9$$

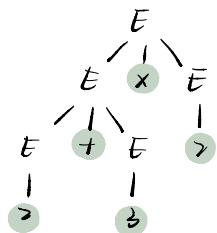
and the derivation

$$E \Rightarrow E+E \Rightarrow 2+E \Rightarrow 2+E \times E \Rightarrow^* 2+3 \times 2$$



We can also have

$$E \Rightarrow E \times E \Rightarrow E+E \times E \Rightarrow^* 2+3 \times 2$$



same word, different parsing trees

The CFA $E \Rightarrow E+E \mid E \times E \mid 0 \mid 1 \mid \dots \mid 9$ is ambiguous as some words such as $2+3 \times 2$ have 2 different parsing trees.

Remark

$$\begin{cases} E \Rightarrow E+E \Rightarrow E+E \times E \Rightarrow 2+E \times E \Rightarrow 2+3 \times E \Rightarrow 2+3 \times 2 \\ E \Rightarrow E+E \Rightarrow 2+E \Rightarrow 2+E \times E \Rightarrow 2+3 \times E \Rightarrow 2+3 \times 2 \end{cases}$$

have identical parsing trees.

DEF A left-most derivation for a word is a derivation where we always apply the rule to the left-most variable of the string.

DEF A CFA is ambiguous if there is a word with 2 different left-most derivations.

Sometimes, we can fix the ambiguity:

$$\begin{array}{lll} E \rightarrow ETTT & & E \rightarrow E+E \\ T \rightarrow T \times F \mid F & = & T \rightarrow EXE \\ F \rightarrow 0111 \cdots 19 & & E \rightarrow 0111 \cdots 19 \\ \uparrow & & \uparrow \\ \text{non ambiguous} & & \text{ambiguous} \end{array}$$

Remark There are CFGs that are inherently ambiguous;
i.e. They cannot be recognized by any non-ambiguous CFG's.

e.g. $\{a^i b^j c^k \mid i=j \text{ or } j=k\}$ see book.

$$\begin{array}{lll} A \rightarrow BC \mid DE \\ C \rightarrow CC \mid \epsilon & D \rightarrow AD \mid \epsilon \\ B \rightarrow aBb \mid \epsilon & E \rightarrow bEc \mid \epsilon \end{array}$$

THM Every regular language is recognized by a CFG.

Since the language is regular, it is described by a regular expression R .

We will use induction on the length of R .

- base case
- ① $R = a$ for some $a \in \Sigma$ $A \rightarrow a$
 - ② $R = \epsilon$ $A \rightarrow \epsilon$
 - ③ $R = \emptyset$

inductive step Assume that every R of length $< m$ can be converted to a CFG.

- ① $R = (R_1)^*$ where $|R_1| < m$
Add $A \rightarrow AA$, $A \rightarrow \epsilon$ to the CFG of R_1 .
- ② $R = R_1 \cup R_2$
Add $\overset{\circ}{S} \rightarrow A_1 \mid A_2$ to the CFG of R_1, R_2
 \downarrow
new start $A_1 \rightarrow \quad A_2 \rightarrow$
- ③ $R = R_1 R_2$
Add $S \rightarrow A_1 A_2$

Chomsky Normal Form

Every CFG can be converted to an equivalent CFG with the start state A s.t. every rule is of the following form:

$B \rightarrow CD$ where $B, C, D \in V$, $C, D \neq A$

$B \rightarrow a$ where $B \in V$, $a \in \Sigma$

$\underset{\text{start}}{A} \rightarrow \epsilon$ only start can go to ϵ

Ex. Which one in Chomsky normal form?

$$\begin{array}{l} A \rightarrow \epsilon \\ A \rightarrow \underline{AB} \times \\ B \rightarrow 0 \\ A \rightarrow 1 \end{array}$$

$$\begin{array}{l} A \rightarrow \epsilon \\ A \rightarrow BC \\ B \rightarrow \underline{0B} \times \\ C \rightarrow 1 \end{array}$$

$$\begin{array}{l} A \rightarrow BC \\ B \rightarrow BC \\ B \rightarrow DB \\ C \rightarrow 1 \\ D \rightarrow 0 \end{array} \quad L = \Sigma^3$$

Convert an arbitrary CFG (V, Σ, R, A) into Chomsky Normal Form:

- ① Add a new start variable A_0 and the rule $A_0 \rightarrow A$
⇒ Guarantee that A_0 is never in the RHS
- ② For every $b \in \Sigma$, add a new variable B and replace all the occurrence of b in the RHS with B . Finally, add $B \rightarrow b$
- ③ If there is a rule $B \rightarrow \epsilon$ where $B \neq A_0$, then we remove it. and for every rule $C \rightarrow \underbrace{u_1 B u_2 B \dots B u_k}_{k \text{ B's}}$, we add all the possible $>^k$ combinations that can be obtained from the RHS by removing some of the B's.
- ④ Consider $B \rightarrow C_1 C_2 \dots C_k$ where $C_1 \dots C_k \in V$ and $k > 2$:

$$B \rightarrow C_1 D_1$$

$$D_1 \rightarrow C_2 D_2$$

$$D_2 \rightarrow C_3 D_3$$

⋮

$$D_{k-2} \rightarrow C_{k-1} C_k$$

where $D_1, D_2 \dots D_{k-2}$ are new variables

- ⑤ $B \rightarrow C$ we remove this and for every $C \rightarrow u$, we add $B \rightarrow u$. We repeat this until we eliminate all such rules.

Push-Down Automaton

Recall: Stack (last in first out)

PDA = NFA + infinite stack

DEF A PDA is a 6-tuple $(Q, \Sigma, \Gamma, \delta, q_0, F)$
 stack alphabet

Everything is similar to NFA except:

- The label on the arrows are of the form

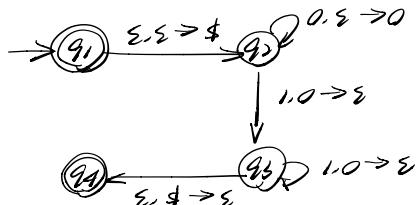
$$\begin{array}{c} Q, A \xrightarrow{\quad} B \xrightarrow{\quad} BEPUSS \\ \text{at } \Sigma \text{ UESS} \quad \text{at } PUSS \end{array}$$

which says if we need a from input, we pop A and push B to stack

$$\delta: Q \times \Sigma_E \times \Gamma_E \rightarrow P(Q \times \Gamma_E) \text{ where } \Sigma_E = \Sigma \cup \{\epsilon\}, \Gamma_E = \Gamma \cup \{\epsilon\}$$

- Acceptance: A PDA accepts a word w if it can read all of w and end at an accept state. (No need to empty stack) non-deterministic

$$\text{Ex. } \Sigma = \{0, 1\} \quad \Gamma = \{0, \$\}$$

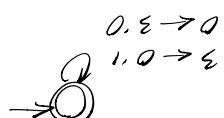


Accept on:

$$\begin{aligned} q_1, \boxed{0} &\xrightarrow{\quad} q_2, \boxed{\$} \xrightarrow{\quad} q_3, \boxed{0} \xrightarrow{\quad} q_4, \boxed{0} \xrightarrow{\quad} q_4, \boxed{\$} \\ &\Rightarrow q_4, \boxed{\$} \xrightarrow{\quad} \text{accept} \end{aligned}$$

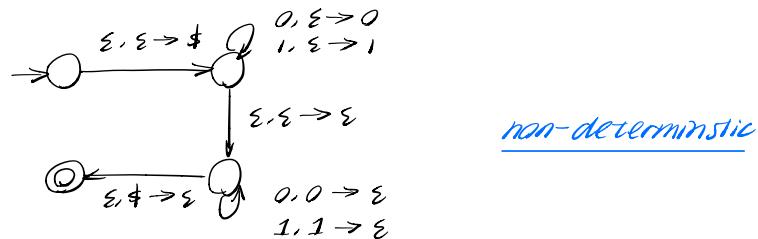
$$L = \{0^n 1^n \mid n \geq 0\}$$

Ex. Design a PDA with one state for the following language L
 {w | no prefix of w contains more 1's than 0's}



Ex. Design a PDA for $L = \{wwr \mid w \in \Sigma^*\}$, $\Sigma = \{0, 1\}$

We use $\$$ to mark the beginning of the stack



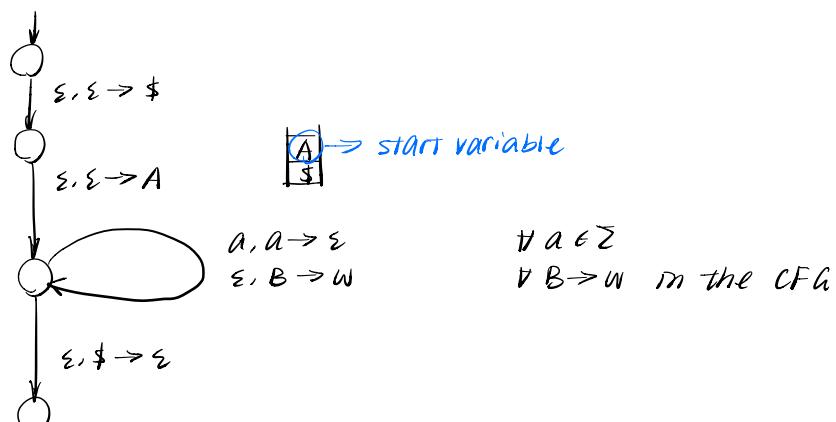
CFA \equiv PDA

THM If L is recognized by a CFA, then there is a PDA for L .

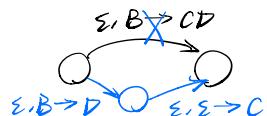
Since L is context free, there is a CFA in Chomsky normal form that generates L
 (V, Σ, R, A)

We start by constructing something that is not quite a PDA: It can pop more than one character at a time

Let $\$$ be a symbol $\$ \notin V \cup \Sigma \cup \{ε\}$

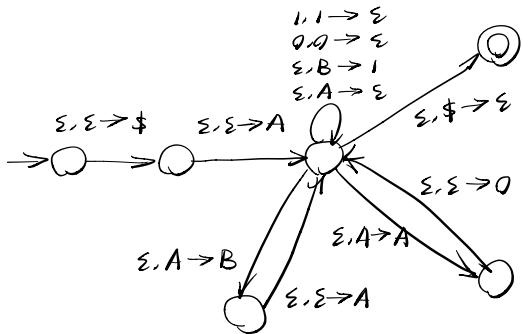


How to fix $B \rightarrow CD$?



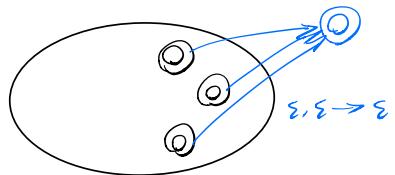
e.g. $A \rightarrow AB \mid 0A \mid \epsilon$

$B \rightarrow 1$



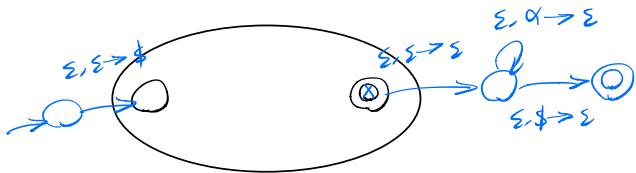
THM If L is recognized by a PDA, then there is a CFA for L .

- ① a single accept state

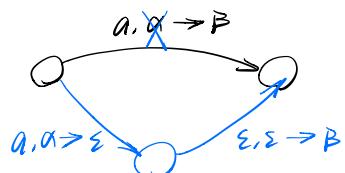


- ② It empties the stack before accepting

PICK a new symbol $\$$



- ③ In each transition $a, a \Rightarrow B$ either $a = \epsilon$ or $B = \epsilon$
(don't push and pop at same time)



Variables For every two states P, Q in PDA, we have a variable $A_{P,Q}$

It is supposed to generate all the words w that can take us from P to Q starting and ending with the empty stack.

$$A_{P,Q} \Rightarrow^* w$$

Start $A_{q_1 \text{ start } q_{accept}}$

\Rightarrow^* the words from start to accept with an empty stack
i.e. Language of the PDA

Rule

- $A_{P,P} \Rightarrow \epsilon$
- $\frac{uw}{A_{P,Q}} \Rightarrow A_{P,u} A_{u,w} A_{w,Q}$



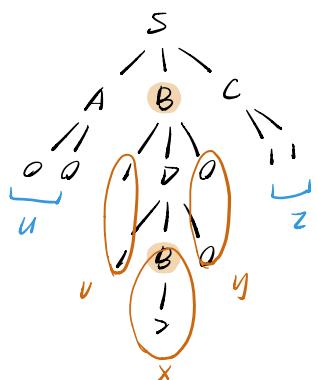
$$A_{P,Q} \Rightarrow a A_{R,S} b$$

Languages that are NOT context-free

Ex. $S \Rightarrow ABC$
 $A \Rightarrow 00$
 $C \Rightarrow 11$

$B \Rightarrow 1D01$
 $D \Rightarrow 1B0$

$$w = 001120011$$



$$\begin{aligned} B &\Rightarrow^* 11B00 \\ B &\Rightarrow^* 111B0000 \\ B &\Rightarrow^* 11111B000000 \end{aligned}$$

Thus, $0011111200000011 \in L$

How long w has to be to guarantee that we will find such R ?

Let $G = (V, \Sigma, R, A)$ for L . Let b be the longest length of the RHS of a rule: $B \Rightarrow \overbrace{\dots}^{\text{length } b}$

Observation 1 Each rule has at most b children

Observation 2 If no such R then the length of the longest path from A to a leaf is at most $|V| + 1$

\Rightarrow at most $b^{|V|}$ leaves

Pumping Lemma of CFG

If L is a CFG, then there \exists a positive p (called pumping constant) such that if $w \in L$ and $|w| \geq p$ then there \exists a decomposition

$$w = uvxyz \quad \text{s.t.}$$

- ① $|vxy| \leq p$
- ② at least one of v and y is not ϵ
- ③ $uv^ixy^iz \in L$

Ex. Show that $L = \{a^n b^n c^n \mid n \geq 0\}$ over $\Sigma = \{a, b, c\}$ is NOT context free.

Suppose it is. Let p be its pumping constant.

Consider $w = a^p b^p c^p \in L$

By Pumping Lemma: $w = uvxyz$ s.t. \dots

Since $|vxy| \leq p$, at least one of the letters a, b, c is missing in vxy .

Then in $uvvxyyz$, at least one of the other two symbols appears more than the missing symbols in xyz . ↗

Ex. Show that both

$$L_1 = \{a^n b^n c^m \mid n, m \geq 0\}$$

$$L_2 = \{a^n b^n c^n \mid n, m \geq 0\}$$

are context-free.

$$\text{for } L_1, \begin{aligned} A &\rightarrow BC \\ B &\rightarrow aBb \mid \epsilon \\ C &\rightarrow cC \mid \epsilon \end{aligned}$$

$$\text{for } L_2, \begin{aligned} A &\rightarrow BC \\ B &\rightarrow bBb \mid \epsilon \\ C &\rightarrow bCc \mid \epsilon \end{aligned}$$

Note. $\{\alpha^n b^n c^n \mid n \geq 0\} = \{\alpha^n b^n c^m \mid n, m \geq 0\} \cap \{\alpha^m b^n c^n \mid n, m \geq 0\}$

NOT CFG

CFG

CFG

\Rightarrow CFG NOT closed under intersection

However, CFG is closed under union and concatenation.

$$L_1 \quad \boxed{A \rightarrow \dots} \quad L_2 \quad \boxed{A' \rightarrow \dots} \quad S \rightarrow AIA' \quad S \rightarrow AA'$$

But NOT closed under complementation

Take a non-context-free language S s.t. $L = L_1 \cap L_2$ where L_1 and L_2 are context-free. $\Rightarrow L^c = (L_1^c \cup L_2^c)$ will be context-free?

Ex. Show that $\{w \in \{0,1\}^* \mid w \in S\}$ is NOT context-free over $\Sigma = \{0,1\}$

NOTE: $\{w \in \{0,1\}^* \mid w \in S\} \neq$

$$w = 0^p 1^p 0^p 1^p \in L$$

$$\underbrace{000 \dots 11 \dots 1}_{\text{first half}} \quad \underbrace{00 \dots 011 \dots 1}_{\text{second half}}$$

① wxy completely in the first half

Then $uvvxyy়z$ is not length $4p + 1v + 1$.

$$\overbrace{00 \dots 11 \dots 1}^{\text{start with } 0} \quad \overbrace{110 \dots 01 \dots 1}^{\text{start with } 1} \\ \overbrace{\text{IVYI}}^{>}$$

② wxy completely in the second half

$$\overbrace{00 \dots 011 \dots 1}^{\text{end with } 0} \quad \overbrace{100 \dots 011 \dots 1}^{\text{end with } 1}$$

③ vxy touches both sides at the first and second half

00 ... 0 11 ... 100 ... 0 11 ... 1

then $i=0 \Rightarrow uxz \in L$

if $uxz = 0^k, 1^k, 0^k, 1^k \Rightarrow k=p$ because the word starts with $\underbrace{0 \dots 0}_p 1$
but $uxz = 0^p, 1^p, 0^p, 1^p = \underline{uvxyz} = w \neq$

Ex. $\{w \mid w \text{ has same number of } a, b, c\}$ over $\Sigma = \{a, b, c\}$ is NOT context free.

$$w = a^p b^p c^p$$