

Properties of Context-Free Languages

- First we will present some closure properties of CF languages and then show some solvable and unsolvable problems related to them.

Closure Properties :

Theorem : If L_1 and L_2 are context-free languages, then so is $L_1 \cup L_2$.

Proof : We need to show there exists a CF grammar that will generate $L_1 \cup L_2$. Let Γ_1 generate L_1 and Γ_2 generate L_2 . Without loss of generality, assume the variables of Γ_1 and L_2 are disjoint. The variables of Γ : $V_1 \cup V_2 \cup \{S\}$ where S is a new variable for Γ — it will be the start symbol for Γ . The productions of Γ are the productions of Γ_1 and Γ_2 plus the productions $S \rightarrow S_1$ and $S \rightarrow S_2$.

Theorem : There exist context-free languages L_1 and L_2 such that $L_1 \cap L_2$ is NOT context-free.

Proof : Because of the word 'exist' in the statement of the theorem, we need to show two CF languages whose intersection is not CF.

This is surprising because as you saw in Lecture 18, regular languages are closed under intersection.

The two languages are : $L_1 = \{a^{[n]} b^{[m]} c^{[m]} \mid n, m > 0\}$, $\Gamma_1 : S \rightarrow Sc, S \rightarrow Xc, X \rightarrow aXb, X \rightarrow ab$
 $L_2 = \{a^{[n]} b^{[n]} c^{[n]} \mid n, m > 0\}$, $\Gamma_2 : S \rightarrow as, S \rightarrow aX, X \rightarrow bXc, X \rightarrow bc$
 $L_1 \cap L_2 = \{a^{[n]} b^{[n]} c^{[n]} \mid n > 0\}$ that we know is not CF (by Bar-Hillel's pumping lemma)

Theorem : There exists a CF language $L \subseteq A^*$ such that $A^* - L$ is NOT context-free.

Proof : Follows from $L_1 \cap L_2 = A^* - ((A^* - L_1) \cup (A^* - L_2))$. If the set subtraction operation yielded a CF language, then $A^* - L_1$ and $A^* - L_2$ would each be CF. Since we already know that the union operation leads to a CF language, that would imply that $(A^* - L_1) \cup (A^* - L_2)$ is CF. But then that would imply that $L_1 \cap L_2$ is CF. Contradiction.

Theorem : If R is a regular language and L is a context-free language, then $R \cap L$ is a context-free language.

Proof : Let T be a positive CF grammar that generates L .

This theorem sounds even more interesting if you remember that every regular language is CF. So, whereas, in general, the intersection of two CF languages is NOT CF, we get a CF intersection if one of the languages is regular.

- Let M be the DFA that accepts R .
- We will construct a positive CF grammar $\tilde{\Gamma}$ that generates $L \cap R$.

$M : A \quad Q \quad \delta, \quad F \quad S$

$\Gamma : A \quad V \quad S \quad \text{productions}$

$\tilde{\Gamma} : A \quad \{p, q\} \quad \tilde{S} \quad \text{productions}$

$\delta \in A \cup V$
and for every ordered pair $p, q \in Q$

- $\tilde{S} \rightarrow S^{q, q}$ for all $q \in F$
- $X^{p, q} \rightarrow \sigma_1^{p, r_1} \sigma_2^{r_1, r_2} \dots \sigma_{n-1}^{r_{n-1}, q}$ for all $r_1, r_2, \dots, r_{n-1}, q \in Q$
- $a^{p, q} \rightarrow a$ for all $a \in A$ and for all $p, q \in Q$ such that $\delta(p, a) = q$

productions
 $X \rightarrow \sigma_1 \sigma_2 \dots \sigma_n$
and Γ and $\tilde{\Gamma}$

About the construction of \tilde{T} on the previous slide, note that it generates a terminal symbol a only when the dfa M contains a state transition for that symbol. This is controlled by the third type of productions specified for \tilde{T} .

- For the proof in the forward direction, we start with a word u that is in $L(\tilde{T})$ and also in R . That is, $u \in L(M)$ and $u \in L(\tilde{T})$. We now want to prove that $u \in L(\tilde{T})$. So we are given $S \xrightarrow[\tilde{T}]{*} u = a_1 a_2 \dots a_n$ with $s(q_1, a_1) = q_2$, $s(q_2, a_2) = q_3, \dots, s(q_n, a_n) = q_{n+1} \in F$. Now we must prove $S \xrightarrow[\tilde{T}]{*} a_1 a_2 \dots a_n$. In terms of the productions specified for \tilde{T} on the previous page, we obviously have

$$\tilde{S} \xrightarrow{\tilde{T}} S^{q_1 q_{n+1}} \xrightarrow[\tilde{T}]{*} a_1^{q_1 q_2} a_2^{q_2 q_3} \dots a_n^{q_n q_{n+1}} \xrightarrow{\tilde{T}} a_1 a_2 \dots a_n$$

This establishes the proof in one direction.

- The proof in the other direction is based on the lemma that when $\sigma \xrightarrow[\tilde{T}]{*} u \in A^*$

then $s^*(p, u) = q$ and when $\sigma \in V$ then we also have $\sigma \xrightarrow[\tilde{T}]{*} u$.

See the text for a proof of the lemma.

because by ③ $a_i^{p,q} \rightarrow a$
when $s(p, a_i) = q$ which is
indeed true by virtue of
④ above.

Sub-Alphabet Erasure :

Let A be the alphabet of a language L . Now consider a subset $P \subseteq A$. We want to delete all of the symbols of P from the words of L . Let $x \in L$. We denote by $Er_p(x)$ what remains of x after erasing from x the symbols of P . We now write

$$Er_p(L) = \{ Er_p(x) \mid x \in L \}$$

Now consider the more specific case of a language generated by a CF grammar T . In this case, we are interested in erasing a subset of terminals. We can prove the following :

If $L \subseteq A^*$ is a CF language and $P \subseteq A$, then $Er_p(L)$ is also CF.

The Word Problem for CF languages :

word problem means figuring out whether $u \in L$ for a $u \in A^*$

- You'll recall that the word problem was unsolvable for a semi-Thue process.
- The word problem was trivially solvable for languages accepted by finite automata (and therefore for languages generated by regular grammars).
- The good news is that the word problem is solvable for context-free languages. This is because of the following theorem:

THEOREM :

Let Δ be a Chomsky normal form grammar with terminals T . For a variable V of Δ let $V \xrightarrow[\Delta]{*} u \in T^*$

Then there is a derivation of u from V in Δ of length $2|u|$.

Proof : We can prove this by induction on $|u|$. When $|u|=1$, then

u is a terminal. So we must have $V \rightarrow a$ for some $a \in T$. The length of this derivation is 2. Now let's consider u with $|u| > 1$. We must have $V \xrightarrow{} Xy \xrightarrow{*} u$. We should therefore expect $X \xrightarrow{\Delta} V$ and $y \xrightarrow{*} w$ with $u = vw$.

By the induction hypothesis, we must have $X = x_1 \Rightarrow x_2 \Rightarrow \dots \Rightarrow x_{2|w|} = v$ and $y = \beta_1 \Rightarrow \beta_2 \Rightarrow \dots \Rightarrow \beta_{2|w|} = w$. Therefore, $V \xrightarrow{} xy \xrightarrow{} x_1 y \xrightarrow{} x_2 y \xrightarrow{} \dots \xrightarrow{} x_{2|w|} y \xrightarrow{} v \beta_1 \xrightarrow{} v \beta_2 \xrightarrow{} \dots \xrightarrow{} v \beta_{2|w|} = vw = u$. But the length of this derivation is $2|w| + 2|v| = 2|u|$.

see Lecture 20 for definition of the length of a derivation

* longest possible path in Δ has $n+2$ nodes.

- $L(T) \neq \emptyset$ if there is a derivation tree T in T of a word $u \in T^*$ such that each path in T contains fewer than $n+2$ nodes. (From Bar-Hillel's pumping lemma) Let u be the shortest possible string in L whose