Finite automata and formal languages (DIT322, TMV028)

Nils Anders Danielsson

2019-03-02

- Closure properties for context-free languages.
- Some algorithms for context-free languages.
- Some undecidable problems.

# Closure properties

- Every regular language is context-free.
- Exercise: Prove this.

#### Assume that

- $\blacktriangleright\ \Sigma_1$  and  $\Sigma_2$  are alphabets and
- $\blacktriangleright \ F \in \Sigma_1 \to \wp(\Sigma_2^*).$

The function F maps symbols to languages. It can be lifted to words and languages:

$$\begin{split} F &\in \Sigma_1^* \to \wp(\Sigma_2^*) \\ F(\varepsilon) &= \{ \varepsilon \} \\ F(aw) &= F(a)F(w) \end{split}$$

$$\begin{array}{l} F\in\wp(\Sigma_1^*)\to\wp(\Sigma_2^*)\\ F(L)=\bigcup_{w\in L}F(w) \end{array}$$

# What is $F(\{ 01 \}^*)$ when $F(0) = \{ a \}$ and $F(1) = \{ b, c \}$ ?

1.  $\{a, b, c\}^*$ 2.  $\{ abc \}^*$ **3**.  $\{ab, ac\}^*$ **4**.  $\{ac, bc\}^*$ 5.  $\{a\}^* \{b,c\}^*$ 6.  $\{a, b\}^* \{c\}^*$ 7.  $\{a\}^* \{bc\}^*$ 8.  $\{ab\}^* \{c\}^*$ 

#### lf

 $\blacktriangleright\ \Sigma_1$  and  $\Sigma_2$  are alphabets,

• 
$$L \subseteq \Sigma_1^*$$
 is context-free,

• 
$$F \in \Sigma_1 \to \wp(\Sigma_2^*)$$
, and

• F(a) is context-free for every  $a \in \Sigma_1$ , then F(L) is context-free.

Idea:

 Replace each terminal a in a grammar for L with the start symbol of a grammar for F(a).

- If L₁ and L₂ are context-free, then L₁ ∪ L₂ is context-free.
- Substitute  $L_i$  for i in  $\{1, 2\}$ .

#### Closure under union

Recall: F(L) is context-free if

- $\blacktriangleright\ \Sigma_1$  and  $\Sigma_2$  are alphabets,
- $L \subseteq \Sigma_1^*$  is context-free,
- $\blacktriangleright \ F \in \Sigma_1 \to \wp(\Sigma_2^*) \text{, and}$
- F(a) is context-free for every  $a \in \Sigma_1$ .

#### In this case:

 ∑<sub>1</sub> = { 1,2 }, ∑<sub>2</sub> is the union of the sets of terminals of some grammars for L<sub>1</sub> and L<sub>2</sub>.

• 
$$L = \{ 1, 2 \} \subseteq \Sigma_1^*$$
 is context-free.

▶ 
$$F(1) = L_1$$
,  $F(2) = L_2$ .

▶ F(1) and F(2) are context-free. Thus  $F(L) = L_1 \cup L_2$  is context-free.

- If L₁ and L₂ are context-free, then L₁L₂ is context-free.
- Substitute  $L_i$  for i in  $\{12\}$ .

- If L is context-free,
  then L\* is context-free.
- Substitute L for 1 in  $\{1\}^*$ .

#### Closure under Kleene plus

- ▶ If L is context-free. then  $L^+$  is context-free.
- Substitute L for 1 in  $\{1\}^+$ .

#### Homomorphisms

Assume that

• 
$$\Sigma_1$$
 and  $\Sigma_2$  are alphabets and

$$\blacktriangleright h \in \Sigma_1 \to \Sigma_2^*.$$

The function h maps symbols to words. It can be lifted to words and languages:

$$\begin{split} & h \in \Sigma_1^* \to \Sigma_2^* & h \in \wp(\Sigma_1^*) \to \wp(\Sigma_2^*) \\ & h(\varepsilon) &= \varepsilon & h(L) = \{ \ h(w) \mid w \in L \} \\ & h(aw) = h(a)h(w) \end{split}$$

The function  $h \in \Sigma_1^* \to \Sigma_2^*$  is a *string homomorphism*.

#### Closure under homomorphism

- If  $L \subseteq \Sigma_1^*$  is context-free, then h(L) is context-free.
- Apply the substitution  $F(a) = \{ h(a) \}$  to L.

Prove that  $L = \{ 01^n 23^n 45^n 6 \mid n \in \mathbb{N} \}$  is not a context-free language over  $\{ 0, 1, 2, 3, 4, 5, 6 \}.$ 

You may use the fact that  $\{ 0^n 1^n 2^n \mid n \in \mathbb{N} \}$  is not a context-free language over  $\{ 0, 1, 2 \}$ .

Hint: Can you find a string homomorphism h for which  $h(L)=\{\ 0^n1^n2^n\mid n\in\mathbb{N}\ \}?$ 

- If  $L_1$  and  $L_2$  are context-free, then  $L_1 \cap L_2$  is *not* necessarily context-free.
- If L<sub>1</sub> and L<sub>2</sub> are context-free, then L<sub>1</sub> \ L<sub>2</sub> is not necessarily context-free.
- If L is a context-free language over Σ, then <u>L</u> = Σ<sup>\*</sup> \ L is *not* necessarily context-free.

- If L is context-free and R is regular, then  $L \cap R$  is context-free.
- If L is context-free and R is regular, then  $L \setminus R$  is context-free.

If  $\Sigma$  is an alphabet,  $R \subseteq \Sigma^*$  is regular and  $L \subseteq \Sigma^*$  is context-free, what can we say about  $R \setminus L$ ?

1. It is always regular.

- 2. It is not necessarily regular, but always context-free.
- 3. It is not necessarily context-free.

Hint:  $\Sigma^* \setminus L = \overline{L}$ .

# Some algorithms

For any context-free language L, given as a context-free grammar  $G=(N,\Sigma,P,S)$ , we can decide if  $L=\emptyset$ :

- A symbol X ∈ N ∪ Σ is generating if X ⇒<sup>\*</sup> w for some w ∈ Σ<sup>\*</sup>.
- $L = \emptyset$  if and only if S is not generating.

### Computing the generating symbols

The set of generating symbols can be computed (perhaps inefficiently) in the following way:

• Let the function  $step \in \wp(N \cup \Sigma) \to \wp(N \cup \Sigma)$  be defined by

$$step(\Gamma) = \left\{ \left. A \right| \begin{array}{l} A \to \alpha \in P, \\ \text{every symbol in } \alpha \text{ is in } \Gamma \end{array} \right\}$$

- Initialise  $\Gamma$  to  $\Sigma$ .
- Repeat until  $step(\Gamma) \subseteq \Gamma$ :
  - Set  $\Gamma$  to  $\Gamma \cup step(\Gamma)$ .
- Return Γ.

Compute the generating symbols of the grammar  $(\{S, A, B\}, \{0, 1\}, P, S)$ , where P is defined in the following way:

 $S \to 0A \mid B$  $A \to 1S \mid \varepsilon$  $B \to AB$ 

S.
 A.
 B.

4. 0.
 5. 1.

For any context-free language L, given as a context-free grammar  $G=(N,\Sigma,P,S)$ , we can decide if  $\varepsilon\in L$ :

- A nonterminal  $A \in N$  is *nullable* if  $A \Rightarrow^* \varepsilon$ .
- We have  $\varepsilon \in L$  if and only if S is nullable.

## Computing the nullable nonterminals

The set of nullable nonterminals can be computed (perhaps inefficiently) in the following way:

• Let the function  $step \in \wp(N) \to \wp(N)$  be defined by

$$step(E) = \left\{ \begin{array}{c} A \mid A \to \alpha \in P, \\ \text{every symbol in } \alpha \text{ is a} \\ \text{nonterminal in } E \end{array} \right\}$$

- Initialise E to  $\emptyset$ .
- Repeat until  $step(E) \subseteq E$ :
  - Set E to  $E \cup step(E)$ .
- Return E.

For any context-free language L, given as a context-free grammar G, and for any nonempty string  $w \in \Sigma^*$ , we can decide if  $w \in L$ .

## The CYK algorithm

- Convert G to a grammar G' = (N, Σ, P, S) in Chomsky normal form such that w ∈ L(G') ⇔ w ∈ L(G).
- Build a CYK table T for G' and w:

• Check if  $S \in T_{1,|w|}$ .

### The CYK algorithm

The table can be computed in the following way:

First set

$$T_{i,i} = \{ A \mid A \rightarrow w_i \in P \}$$
  
for each  $i \in \{1, ..., |w|\}.$   
 $\blacktriangleright$  Then set

$$T_{i,j} = \left\{ \begin{array}{l} A \mid k \in \left\{ \begin{array}{l} i, ..., j - 1 \end{array} \right\}, \\ B \in T_{i,k}, C \in T_{k+1,j}, \\ A \to BC \in P \end{array} \right\}$$

for all  $i,j\in\{1,...,|w|\}$  satisfying j-i+1=2.

• Repeat the previous step for j - i + 1 = 3, 4 and so on up to |w|.



#### An example of dynamic programming.

Consider the following CYK table:

$$\begin{cases} S \\ \emptyset & \{T \\ \emptyset & \{S \} & \emptyset \\ \{T,Z \} & \{U,O \} & \{U,O \} & \{T,Z \} \\ \hline 0 & 1 & 1 & 0 \\ \end{cases}$$

Construct a parse tree for the string 0110, given the information that at least the following productions exist in the grammar:  $S \rightarrow ZT, S \rightarrow OU, T \rightarrow SZ$ .

- A potential problem:
  The size of G' can be quadratic in the size of G.
- A variant of the algorithm that does not use the UNIT transformation can be devised:
  - Time complexity:  $O(|G||w|^3)$ .
  - Space complexity:  $O(|G||w|^2)$ .

See Lange and Leiß.

# Some undecidable problems

The following things cannot, in general, be determined (using, say, a Turing machine):

- If a context-free grammar is ambiguous.
- If a context-free language, given by a context-free grammar, is *inherently* ambiguous.
- ▶ If L(G<sub>1</sub>) = L(G<sub>2</sub>) for two context-free grammars G<sub>1</sub> and G<sub>2</sub>.

• ...

If you want to know more about why certain problems are undecidable, then you might be interested in the course *Computability*.

- Closure properties for context-free languages.
- Some algorithms for context-free languages.
- Some undecidable problems.



- Pushdown automata.
- Turing machines.