

Finite automata and formal languages (DIT322, TMV028)

Nils Anders Danielsson

2020-02-20

Today

Context-free grammars: syntax and semantics.

Context-free grammars

Context-free grammars (CFGs)

- ▶ The context-free languages are those that can be described by CFGs.
- ▶ Every regular language is context-free.
- ▶ Some context-free languages are not regular.
- ▶ CFGs are for instance used to specify the syntax of some programming languages.
 - ▶ One example: Haskell.
- ▶ Parser generators often use (restricted) CFGs.

Syntax

Context-free grammars

A context-free grammar has the form (N, Σ, P, S) :

- ▶ N is a finite set of *nonterminals*.
- ▶ Σ is a finite set of *terminals* satisfying
 $\Sigma \cap N = \emptyset$.
- ▶ $P \subseteq N \times (N \cup \Sigma)^*$ is a finite set of
productions.
- ▶ The start symbol $S \in N$.

Notation

- ▶ A production (A, α) can be written $A \rightarrow \alpha$.
- ▶ Multiple productions $A \rightarrow \alpha_1, \dots, A \rightarrow \alpha_n$ can be written $A \rightarrow \alpha_1 \mid \dots \mid \alpha_n$ (if $n \geq 2$).

Which of the following expressions are well-formed context-free grammars?

1. $(\mathbb{N}, \{ a, b \}, P, 0)$, where P contains the following productions: $0 \rightarrow a1, 1 \rightarrow b$.
2. $(\{ 0, 1 \}, \{ a, b \}, P, 0)$, where P contains the following productions: $0 \rightarrow a1, 1 \rightarrow b$.
3. $(\{ 0, 1 \}, \{ 0, 1 \}, P, 0)$, where P contains the following productions: $0 \rightarrow 01, 1 \rightarrow 1$.
4. $(\{ 0, 1 \}, \{ 0', 1' \}, P, 0)$, where P contains the following productions: $0 \rightarrow 01, 1 \rightarrow 1 \mid 0$.
5. $(\{ 0, 1 \}, \{ 0', 1' \}, P, 2)$, where P contains the following productions: $0 \rightarrow 01, 1 \rightarrow 1 \mid 0$.

Examples

An example

A context-free grammar for the non-regular language $\{ 0^n 1^n \mid n \in \mathbb{N} \}$ over $\{ 0, 1 \}$:

$$(\{ S \}, \{ 0, 1 \}, S \rightarrow 0S1 \mid \varepsilon, S)$$

An example

A context-free grammar for the non-regular language $\{ 0^n 1^n \mid n \in \mathbb{N} \}$ over $\{ 0, 1 \}$:

$$(\{ S \}, \{ 0, 1 \}, S \rightarrow 0S1 \mid \varepsilon, S)$$

Generated strings:

- ▶ $\varepsilon.$
- ▶ $0\varepsilon1 = 01.$
- ▶ $0011.$
- ▶ \vdots

An example

A context-free grammar for the non-regular language $\{ 0^n 1^n \mid n \in \mathbb{N} \}$ over $\{ 0, 1 \}$:

$$(\{ S \}, \{ 0, 1 \}, S \rightarrow 0S1 \mid \varepsilon, S)$$

An inductive definition of the language $L \subseteq \{ 0, 1 \}^*$ generated by the grammar:

$$\frac{w \in L}{0w1 \in L} \qquad \qquad \overline{\varepsilon \in L}$$

Another example

Consider the grammar $(\{ S, A \}, \{ 0, 1 \}, P, S)$,
where P is defined in the following way:

$$S \rightarrow 0A1 \mid \varepsilon$$

$$A \rightarrow 1A0 \mid S$$

Another example

Consider the grammar $(\{ S, A \}, \{ 0, 1 \}, P, S)$,
where P is defined in the following way:

$$S \rightarrow 0A1 \mid \varepsilon$$

$$A \rightarrow 1A0 \mid S$$

Sentential forms:

- ▶ $S.$
- ▶ $\varepsilon.$
- ▶ $0A1.$
- ▶ $01A01.$
- ▶ $01S01.$
- ▶ $0101.$
- ▶ \vdots

Another example

Consider the grammar $(\{ S, A \}, \{ 0, 1 \}, P, S)$,
where P is defined in the following way:

$$S \rightarrow 0A1 \mid \varepsilon$$

$$A \rightarrow 1A0 \mid S$$

An inductive definition of the languages
 $L_S, L_A \subseteq \{ 0, 1 \}^*$ generated by S and A :

$$\frac{w \in L_A}{0w1 \in L_S} \qquad \qquad \frac{}{\varepsilon \in L_S}$$

$$\frac{w \in L_A}{1w0 \in L_A} \qquad \qquad \frac{w \in L_S}{w \in L_A}$$

Construct a context-free grammar for the language $\{ 0^{3n}1^{2n} \mid n \in \mathbb{N} \}$ over $\{ 0, 1 \}$ by filling in the missing part of the following definition.

$(\{ S \}, \{ 0, 1 \}, S \rightarrow ???, S)$

Semantics

Derivations

For the grammar $G = (N, \Sigma, P, S)$ one can define the following two binary relations on $(N \cup \Sigma)^*$ inductively:

$$\frac{\alpha, \beta \in (N \cup \Sigma)^* \quad A \in N \quad (A, \gamma) \in P}{\alpha A \beta \Rightarrow \alpha \gamma \beta}$$

$$\frac{}{\alpha \Rightarrow^* \alpha} \qquad \frac{\alpha \Rightarrow \beta \quad \beta \Rightarrow^* \gamma}{\alpha \Rightarrow^* \gamma}$$

The language $L(G) = \{ w \in \Sigma^* \mid S \Rightarrow^* w \}.$

Leftmost derivations

A variant:

$$\frac{w \in \Sigma^* \quad A \in N \quad \alpha \in (N \cup \Sigma)^* \\ (A, \beta) \in P}{wA\alpha \Rightarrow_{\text{Im}} w\beta\alpha}$$

$$\frac{}{\alpha \Rightarrow_{\text{Im}}^* \alpha} \qquad \frac{\alpha \Rightarrow_{\text{Im}} \beta \quad \beta \Rightarrow_{\text{Im}}^* \gamma}{\alpha \Rightarrow_{\text{Im}}^* \gamma}$$

Which of the following propositions are valid?

1. $A \Rightarrow^* \beta \Leftrightarrow A \Rightarrow_{\text{Im}}^* \beta$
2. $A \Rightarrow^* w \Leftrightarrow A \Rightarrow_{\text{Im}}^* w$
3. $\alpha \Rightarrow^* \beta \Leftrightarrow \beta \Rightarrow^* \alpha$
4. $\exists w \in \Sigma^*. \varepsilon \Rightarrow w$
5. $\exists w \in \Sigma^*. \varepsilon \Rightarrow_{\text{Im}}^* w$

A bug

The course text book states that

$$A \Rightarrow^* \beta \quad \Leftrightarrow \quad A \Rightarrow_{\text{Im}}^* \beta$$

holds. Do not trust everything that you read.

Recursive inference

If the grammar $G = (N, \Sigma, P, S)$, then one can define certain languages over Σ inductively:

- ▶ The language generated by the nonterminal $A \in N$, $L(G, A)$.
- ▶ The language generated by a list $\alpha \in (N \cup \Sigma)^*$, $L_L(G, \alpha)$.

Recursive inference

$$\frac{(A, \alpha) \in P \quad w \in L_L(G, \alpha)}{w \in L(G, A)}$$

$$\frac{}{\varepsilon \in L_L(G, \varepsilon)} \qquad \frac{a \in \Sigma \quad w \in L_L(G, \alpha)}{aw \in L_L(G, a\alpha)}$$

$$\frac{A \in N \quad v \in L(G, A) \quad w \in L_L(G, \alpha)}{vw \in L_L(G, A\alpha)}$$

Recursive inference

Consider the grammar $(\{ A, B \}, \{ 0, 1 \}, P, A)$,
where P is defined in the following way:

$$A \rightarrow 0B0$$

$$B \rightarrow 1A1 \mid \varepsilon$$

Recursive inference

Recall:

$$P = \{ A \rightarrow 0B0, B \rightarrow 1A1, B \rightarrow \varepsilon \}$$

A derivation:

$$\frac{\frac{\frac{\frac{B \rightarrow \varepsilon \in P}{\varepsilon \in L_L(G, \varepsilon)} \quad \frac{\varepsilon \in L_L(G, \varepsilon)}{0 \in L_L(G, 0)}}{\varepsilon \in L(G, B)} \quad \frac{0 \in L_L(G, 0)}{00 \in L_L(G, 0B0)}}{A \rightarrow 0B0 \in P} \quad \frac{00 \in L_L(G, 0B0)}{00 \in L(G, A)}$$

Due to lack of space I have omitted
“ $a \in \Sigma$ ” and “ $A \in N$ ”.

Which of the following propositions are true?

- 1. $1001 \in L(G, A)$
- 3. $00 \in L_L(G, AB)$
- 2. $1001 \in L(G, B)$
- 4. $0000 \in L_L(G, AB)$

Hint: Try to construct derivation trees.

Parse trees

Consider the following definitions again:

$$\frac{(A, \alpha) \in P \quad w \in L_L(G, \alpha)}{w \in L(G, A)}$$

$$\frac{}{\varepsilon \in L_L(G, \varepsilon)} \qquad \frac{a \in \Sigma \quad w \in L_L(G, \alpha)}{aw \in L_L(G, a\alpha)}$$

$$\frac{A \in N \quad v \in L(G, A) \quad w \in L_L(G, \alpha)}{vw \in L_L(G, A\alpha)}$$

Parse trees

Parse trees:

$$\frac{(A, \alpha) \in P \quad ts \in P_L(G, \alpha)}{\text{node}(A, ts) \in P(G, A)}$$

$$\frac{}{\text{nil} \in P_L(G, \varepsilon)} \qquad \frac{a \in \Sigma \quad ts \in P_L(G, \alpha)}{\text{term}(a, ts) \in P_L(G, a\alpha)}$$

$$\frac{A \in N \quad t \in P(G, A) \quad ts \in P_L(G, \alpha)}{\text{nonterm}(t, ts) \in P_L(G, A\alpha)}$$

Recursive inference

Recall:

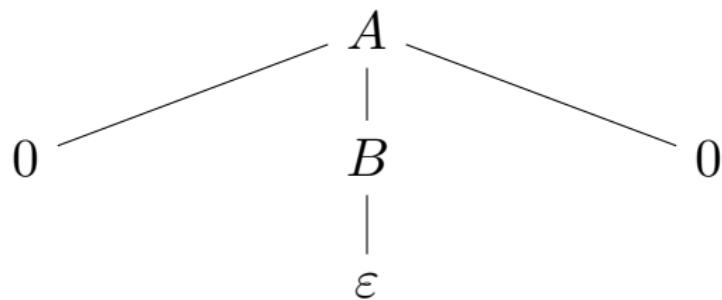
$$\frac{\overline{B \rightarrow \varepsilon \in P} \quad \overline{\varepsilon \in L_L(G, \varepsilon)} \quad \overline{\varepsilon \in L_L(G, \varepsilon)}}{\underline{\varepsilon \in L(G, B)} \quad \overline{0 \in L_L(G, 0)}}$$
$$\frac{\overline{A \rightarrow 0B0 \in P} \quad \overline{0 \in L_L(G, B0)}}{\underline{00 \in L_L(G, 0B0)}}$$
$$00 \in L(G, A)$$

A corresponding parse tree in $P(G, A)$:

`node(A, term(0, nonterm(node(B, nil), term(0, nil)))))`

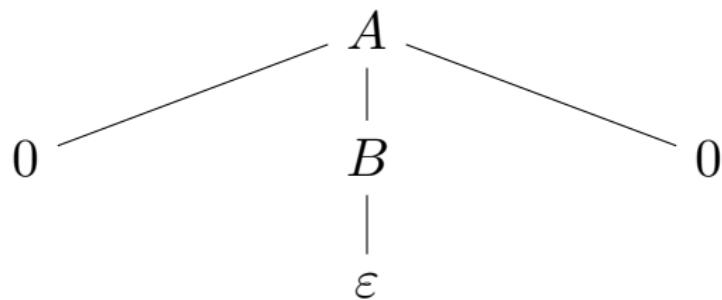
Recursive inference

A different way to present the parse tree:



Recursive inference

A different way to present the parse tree:



This parse tree *yields* the string 00.

Parse trees

The yield of a parse tree:

$$yield \in P(G, A) \rightarrow \Sigma^*$$

$$yield(\text{node}(A, ts)) = yield_L(ts)$$

$$yield_L \in P_L(G,) \rightarrow \Sigma^*$$

$$yield_L(\text{nil}) = \varepsilon$$

$$yield_L(\text{term}(a, ts)) = a \ yield_L(ts)$$

$$yield_L(\text{nonterm}(t, ts)) = yield(t) \ yield_L(ts)$$

Yields containing nonterminals

The inductive definitions of recursive inference and parse trees can be extended to support strings containing both terminals and nonterminals:

$$\overline{A \in L_N(G, A)} \qquad \overline{\text{leaf}(A) \in P_N(G, A)}$$

$$yield \in P_N(G, A) \rightarrow (N \cup \Sigma)^*$$

$$yield(\text{leaf}(A)) = A$$

$$yield(\text{node}(A, ts)) = yield_L(ts)$$

$$yield_L \in P_{NL}(G,) \rightarrow (N \cup \Sigma)^*$$

:

Which of the following propositions are valid?

1. $\forall t \in P(G, A). \text{yield}(t) \in L(G, A)$
2. $\forall ts \in P_L(G, \alpha). \text{yield}_L(ts) \in L_L(G, \alpha)$
3. $\forall \alpha \in (N \cup \Sigma)^*. A \Rightarrow^* \alpha \Leftrightarrow \alpha \in L(G, A)$
4. $\forall \alpha \in (N \cup \Sigma)^*. A \Rightarrow^* \alpha \Leftrightarrow \alpha \in L_N(G, A)$
5. $w \in L_L(G, \alpha\beta) \Leftrightarrow$
 $\exists u \in L_L(G, \alpha), v \in L_L(G, \beta). w = uv$
6. $uv \in L_L(G, \alpha\beta) \Leftrightarrow$
 $u \in L_L(G, \alpha) \wedge v \in L_L(G, \beta)$

Proofs about
grammars

A proof

Recall:

$$G = (\{ S \}, \{ 0, 1 \}, S \rightarrow 0S1 \mid \varepsilon, S)$$

$$\frac{w \in L}{\overline{0w1 \in L}} \qquad \qquad \qquad \overline{\varepsilon \in L}$$

Let us prove that $L(G, S) \subseteq L$.

A proof

Let us prove $\forall w \in L(G, S)$. $w \in L$ by complete induction on the length of the string:

- ▶ $w \in L(G, S)$ implies that
 $w \in L_L(G, \varepsilon)$ or $w \in L_L(G, 0S1)$.
- ▶ If $w \in L_L(G, \varepsilon)$, then $w = \varepsilon \in L$.
- ▶ If $w \in L_L(G, 0S1)$, then...

A proof

- ▶ If $w \in L_L(G, 0S1)$, then:
 - ▶ $w = 0w'$ for some $w' \in L_L(G, S1)$.
 - ▶ $w' = uv$ for some
 $u \in L(G, S), v \in L_L(G, 1)$.
 - ▶ $v = 1$.
 - ▶ $|u| < |w|$, so by the inductive hypothesis
 $u \in L$.
 - ▶ Thus $w = 0u1 \in L$.

A proof

- ▶ Another kind of induction can also be used: induction on the structure of the recursive inference.
- ▶ Exercise (optional, hard):
 - ▶ Write down a formula for this kind of induction.
 - ▶ Use this kind of induction to prove $L(G, S) \subseteq L$.

Prove $L \subseteq L(G, S)$,

i.e. $\forall w \in L. w \in L(G, S)$,

by induction on the structure of L .

$$G = (\{ S \}, \{ 0, 1 \}, S \rightarrow 0S1 \mid \varepsilon, S)$$

$$\frac{w \in L}{0w1 \in L} \qquad \qquad \frac{}{\varepsilon \in L}$$

Today

- ▶ Context-free grammars.
- ▶ Derivations.
- ▶ Left-most derivations.
- ▶ Recursive inference.
- ▶ Parse trees.
- ▶ Proofs about grammars.

Next lecture

- ▶ More about context-free grammars.