

Formal Languages & Compiler Design

Homework 1

Stefan Stefanache

October 26, 2022

1. Given the grammar

$$G = (\{S, H\}, \{b, c, d, e\}, \{S \rightarrow b^2Se \mid H, H \rightarrow cHd^2 \mid cd\}, S) \quad (1)$$

find the language it generates.

Proof. Let us define the language

$$L = \{b^{2n}c^{m+1}d^{2m+1}e^n \mid n, m \in \mathbb{N}\} = \{L_{n,m} \mid n, m \in \mathbb{N}\} \quad (2)$$

and label G 's productions as follows:

$$S \rightarrow b^2Se \quad (3)$$

$$S \rightarrow H \quad (4)$$

$$H \rightarrow cHd^2 \quad (5)$$

$$H \rightarrow cd \quad (6)$$

Our goal is proving that $L = L(G)$ using the double inclusion technique.

Step 1: Show that $L \subseteq L(G)$:

Using induction, one can easily show that for all $n \in \mathbb{N}$,

$$H \xrightarrow[(5)]{n} c^n H d^{2n} \xrightarrow[(6)]{1} c^{n+1} d^{2n+1} \quad (7)$$

$$S \xrightarrow[(3)]{n} b^{2n} S e^n \xrightarrow[(4)]{1} b^{2n} H e^n \quad (8)$$

Therefore, we have that

$$S \xrightarrow[(8)]{*} b^{2k} H e^k \xrightarrow[(7)]{*} b^{2k} c^{p+1} d^{2p+1} e^k, \quad \forall p, k \in \mathbb{N}$$

As a result,

$$S \xRightarrow{*} w, \quad \forall w \in L$$

so

$$L \subseteq L(G) \tag{9}$$

Step 2: Show that $L \supseteq L(G)$:

All elements in $L(G)$ are words/sequences derived from S . Therefore, we have to show that all elements that can be derived from S and don't contain any nonterminal symbols are contained by L .

We start by studying how the elements derived from H look like. After applying the (5) rule for $m \in \mathbb{N}$ times, the result is of the form $c^m H d^{2m}$. To get rid of H , we apply (6) one time to get a word/sequence of the form $c^{m+1} d^{2m+1}$. This is the only kind of word/sequence that can be obtained starting from H .

Now, we continue by looking at the possible elements derived from S . Analogously, we apply rule (3) for $n \in \mathbb{N}$ times to get to the form $b^{2n} S e^n$. At this point, our only move is to transform the S into H using (4). Finally, we transform H like previously discussed, to obtain a final word/sequence of the form $b^{2n} c^{m+1} d^{2m+1} e^n \in L$, with $n, m \in \mathbb{N}$.

Since all possible words/sequences derived from S can be found in L , we have that

$$L \supseteq L(G) \tag{10}$$

Finally, from (9) and (10), we have that $L = L(G)$, so the language generated by the grammar G is given by (2). \square

2. Find grammars that generate the following languages:

- A.** $L_1 = \{x^n y^n \mid n \in \mathbb{N}\}$, with proof
- B.** $L_2 = \{a^n b^{2n} \mid n \in \mathbb{N}\}$, with proof
- C.** $L_3 = \{a^n b^m \mid n, m \in \mathbb{N}^*\}$, with proof using regular grammar
- D.** $L_4 = \{x^{2n} \mid n \in \mathbb{N}\}$, $L'_4 = \{x^{2n} \mid n \in \mathbb{N}^*\}$, with proof using regular grammar
- E.** \mathbb{N}
- F.** All arithmetic expressions containing a as operand, $+$, $*$ as operators and $()$.

Proof. For all proofs, we'll be using a similar method as in the previous exercise: find a grammar G_1 and then prove that $L_1 = L(G_1)$ using the double inclusion technique.

A. Let us define the grammar

$$G_1 = (\{A\}, \{x, y\}, \{A \rightarrow xAy \mid \epsilon\}, A) \quad (11)$$

and label its rules as

$$A \rightarrow xAy \quad (12)$$

$$A \rightarrow \epsilon \quad (13)$$

Step 1. Show that $L_1 \subseteq L(G_1)$:

Using induction, one can easily show that $A \xRightarrow[(12)]{n} x^n A y^n$.

Since $x^n A y^n \xRightarrow[(13)]{1} x^n y^n$, we have that

$$A \xRightarrow{*} x^n y^n \in L_1 \quad (14)$$

Therefore, $A \xRightarrow{*} w, \forall w \in L_1$, so

$$L_1 \subseteq L(G_1) \quad (15)$$

Step 2. Show that $L_1 \supseteq L(G_1)$:

We have to prove that each element derived from A that doesn't contain any nonterminal symbols is contained by L_1 . Starting from A , we apply rule (12) for $n \in \mathbb{N}$ times and reach a result of the form $x^n A y^n$. To get rid of the nonterminal symbol A , we use rule (13) to obtain the word/sequence $x^n y^n$.

Since this is the only possible way of sequencing transformations (note that this also works for $n = 0$), all elements derived from A with no nonterminal terms are of the form $x^n y^n \in L_1$. Therefore,

$$L_1 \supseteq L(G_1) \quad (16)$$

Using (15) and (16), we have that $L_1 = L(G_1)$, so we've proved that the language L_1 is generated by the grammar given by (11).

B. Let us define the grammar

$$G_2 = (\{B\}, \{a, b\}, \{B \rightarrow aBb^2 \mid ab^2\}, B) \quad (17)$$

and label its rules as

$$B \rightarrow aBb^2 \quad (18)$$

$$B \rightarrow ab^2 \quad (19)$$

Step 1. Show that $L_2 \subseteq L(G_2)$:

Let $n \in \mathbb{N}^*$. Using induction, one can easily show that

$$B \xrightarrow[n-1]{(18)} a^{n-1}Bb^{2n-2}$$

Since $a^{n-1}Bb^{2n-2} \xrightarrow[(19)]{1} a^n b^{2n}$, we have that

$$B \xrightarrow{*} a^n b^{2n} \in L_2 \quad (20)$$

Therefore, $B \xrightarrow{*} w, \forall w \in L_2$, so

$$L_2 \subseteq L(G_2) \quad (21)$$

Step 2. Show that $L_2 \supseteq L(G_2)$:

We have to prove that each element derived from B that doesn't contain any nonterminal symbols is contained by L_2 . Starting from B , we apply rule (18) for $n-1 \in \mathbb{N}$ times and reach a result of the form $a^{n-1}Bb^{2n-2}$. To get rid of the nonterminal symbol B , we use rule (19) to obtain the word/sequence $a^n b^{2n}$.

Since this is the only possible way of sequencing transformations (note that this also works for $n=0$), all elements derived from B with no nonterminal terms are of the form $a^n b^{2n} \in L_2$. Therefore,

$$L_2 \supseteq L(G_2) \quad (22)$$

Using (21) and (22), we have that $L_2 = L(G_2)$, so we've proved that the language L_2 is generated by the grammar given by (17).

C. Let us define the **regular** grammar

$$G_3 = (\{A, B\}, \{a, b\}, \{A \rightarrow aA \mid aB, B \rightarrow bB \mid b\}, A) \quad (23)$$

and label its rules as

$$A \rightarrow aA \quad (24)$$

$$A \rightarrow aB \quad (25)$$

$$B \rightarrow bB \quad (26)$$

$$B \rightarrow b \quad (27)$$

Step 1. Show that $L_3 \subseteq L(G_3)$:

Let $n, m \in \mathbb{N}^*$. Using induction twice, one can easily show that

$$A \xrightarrow[(24)]{n-1} a^{n-1}A \xrightarrow[(25)]{1} a^nB \xrightarrow[(26)]{m-1} a^nb^{m-1}B \xrightarrow[(27)]{1} a^nb^m$$

Therefore, we have that

$$A \xRightarrow{*} a^nb^m \in L_3 \quad (28)$$

Therefore, $A \xRightarrow{*} w, \forall w \in L_3$, so

$$L_3 \subseteq L(G_3) \quad (29)$$

Step 2. Show that $L_3 \supseteq L(G_3)$:

We have to prove that each element derived from A that doesn't contain any nonterminal symbols is contained by L_3 . Starting from A , we apply rule (24) for $n - 1 \in \mathbb{N}$ times and reach a result of the form $a^{n-1}A$. To get rid of the nonterminal symbol A , we use rule (19) to obtain the new form a^nB .

Now, we can apply rule (26) for $m - 1 \in \mathbb{N}$ times to reach the form $a^nb^{m-1}B$. Getting rid of the nonterminal B is done by using the rule (27) and reaching the final form $a^nb^m \in L_3$.

Since this is the only possible way of sequencing transformations, all elements derived from A with no nonterminal terms are of the form $a^nb^m \in L_3$. Therefore,

$$L_3 \supseteq L(G_3) \quad (30)$$

Using (29) and (30), we have that $L_3 = L(G_3)$, so we've proved that the language L_3 is generated by the grammar given by (23).

D. Let us define the **regular** grammar

$$G_4 = (\{S, A, B\}, \{x\}, \{S \rightarrow xA \mid \epsilon, A \rightarrow xB \mid x, B \rightarrow xA\}, S) \quad (31)$$

and label its rules as

$$S \rightarrow xA \quad (32)$$

$$S \rightarrow \epsilon \quad (33)$$

$$A \rightarrow xB \quad (34)$$

$$A \rightarrow x \quad (35)$$

$$B \rightarrow xA \quad (36)$$

One trick we're using here is making sure that A is always accompanied by an odd power of x . That way, we either use (35) to transform into something of the form x^{2k} , or transform it into xB using (34) and cycle back with (36), getting something of the form $x^{2k+1}A$ again. By doing it this way, we make sure that only odd powers are covered.

Step 1. Show that $L_4 \subseteq L(G_4)$:

To enable the cycle trick discussed above, we'll prove by induction that

$$A \xRightarrow{*} x^{2n+1}, \forall n \in \mathbb{N} \quad (37)$$

The base case for $n = 0$ holds since $A \xRightarrow[(32)]{1} x$. Now, let's take $k \in \mathbb{N}$ and assume that

$$A \xRightarrow{*} x^{2k+1} \quad (38)$$

Then,

$$A \xRightarrow[(34)]{1} xB \xRightarrow[(36)]{1} x^2A \xRightarrow[(38)]{*} x^{2k+3}$$

This proves that $P(k) \implies P(k+1)$ is true, where

$$P(n) : A \xRightarrow{*} x^{2n+1}$$

Therefore, since both $P(0)$ and $P(k) \implies P(k+1)$ hold we've proved that (37) is true. Now, we have that

$$S \xRightarrow[(32)]{1} xA \xRightarrow[(37)]{*} x^{2n} \in L_4$$

Therefore, $S \xRightarrow{*} w, \forall w \in L_4$, so

$$L_4 \subseteq L(G_4) \quad (39)$$

Step 2. Show that $L_4 \supseteq L(G_4)$:

We have to prove that each element derived from S that doesn't contain any nonterminal symbols is contained by L_4 . Starting from S , we apply rule (32) to obtain xA . Now, we sequentially apply the rules (34) and (36) for $k \in \mathbb{N}$ times to obtain something of the form $x^{2k+1}A$. Now, we can apply (35) to get rid of the nonterminal symbol A and get the final sequence of the form $x^{2k} \in L_4$.

We can also start from S and apply (33) to obtain $\epsilon = x^0 \in L_4$.

Since these are the only possible way of sequencing transformations, all elements derived from A with no nonterminal terms are of the form $a^n b^m \in L_3$. Therefore,

$$L_4 \supseteq L(G_4) \quad (40)$$

Using (39) and (40), we have that $L_4 = L(G_4)$, so we've proved that the language L_4 is generated by the grammar given by (31).

We do exactly the same thing for L'_4 , but by removing the (33) rule from the production set of the grammar, since we're now dealing with only positive powers of x . One can define the adapted **regular** grammar

$$G'_4 = (\{S, A, B\}, \{x\}, \{S \rightarrow xA, A \rightarrow xB \mid x, B \rightarrow xA\}, S) \quad (41)$$

and follow the same steps as above to prove that $L'_4 = L(G'_4)$.

E. Let us define the grammar

$$G_5 = (\{S, A, B, C, D\}, \Sigma, P, S) \quad (42)$$

where the nonterminal symbols are the 10 digits,

$$\Sigma = \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\} \quad (43)$$

and the production set is given by

$$P = \{S \rightarrow 0 \mid A, A \rightarrow B \mid BC, C \rightarrow D \mid DC, D \rightarrow 0 \mid B, B \rightarrow 1 \mid 2 \mid 3 \mid \dots \mid 9\} \quad (44)$$

This grammar generates the set of natural numbers \mathbb{N} ($L(G_5) = \mathbb{N}$) and its rules can be described as follows: B denotes a non-zero digit, D denotes any digit, C denotes a sequence of digits, and A denotes a natural positive number. S can be transformed to a positive number (A) or 0.

F. Let us define the grammar

$$G_6 = (\{S, A, b\}, \{a, +, *, (,)\}, P, S) \quad (45)$$

Using the production set

$$P = \{S \rightarrow S + a \mid A, A \rightarrow A * a \mid B, B \rightarrow (S) \mid a\} \quad (46)$$

we have that $L(G_6)$ covers all arithmetic expressions containing a as an operand, $+$, $*$ as operators and $()$. Compound statements can be enabled by the cyclical nature of the rules $S \rightarrow A \rightarrow B \rightarrow S \rightarrow \dots$. The first rule enables sums, the second products, and the third compounds statements.

□