

Decidable Problems Concerning Context-Free Languages

Topics

- **Problem 1:** describe algorithms to test whether a CFG generates a particular string
- **Problem 2** describe algorithms to test whether the language generated by a CFG is empty.
- **Problem 3:** describe algorithms to test whether an arbitrary string is an element of a context free language, (i.e., is there a CFG G such that a string $w \in L(G)$?)
- **Problem 4:** for two CFLs L_1 and L_2 is $L_1 = L_2$ true?

Note

Problem 1 differs from Problem 3 because in Problem 3 the language is given while in Problem 1 the grammar is given.

Problem 1 string generation problem

Problem: For a given CFG grammar $G = (V, \Sigma, R, S)$ and string $w \in \Sigma^*$, does G generate w (i.e., is $S \xRightarrow{*} w$ true?)

Language: $A_{CFG} = \{\langle G, w \rangle \mid G \text{ is a CFG that generates string } w\}$

Theorem 4.7 (4.6)

A_{CFG} is a decidable language

Proof ideas: For a CFG G and a string w :

- **First idea:** Go through all derivations generated by G checking whether any is a derivation of w .

Since there are infinitely many derivations this idea does not work. If G does not generate w the algorithm doesn't halt. I.e, this idea provides a recognizer but not a decider

A better idea

- Make the recognizer a decider. For that we need to ensure that the algorithm tries only finitely many derivations.

Fact 1

If G is a CFG in Chomsky normal form then for any $w \in L(G)$ where $|w| = n$ exactly $2n - 1$ steps are required for any derivation of w

Proof:

1. Derivation rules of a Chomsky normal form are of the form:
 $A \rightarrow A_1A_2 \mid A \rightarrow a$.
2. The rule $A \rightarrow A_1A_2$ adds 1 to the length of w . That is, if $|w| = n$ then $S \Rightarrow A_1A_2 \Rightarrow A_1A_2A_3 \Rightarrow A_1A_2 \dots A_n$, using $n - 1$ steps.
3. To eliminate A_1, A_2, \dots, A_n by rules of the form $A \rightarrow a$ we need another n steps.

Conclusion: only $2n - 1$ steps are required.

Checking CFG derivations

- Convert G into Chomsky normal form
- For a string w of length $|w| = n$, $n > 1$, check all derivations with $2n-1$ steps to determine whether G generates w
- For a string w , of length $|w| \leq 1$, check all one-step derivations to determine whether G generates w .

Note: since we can convert G into a Chomsky normal form (see Section 2.1), this is a good idea

Proof of Theorem 4.7

The TM S that decide A_{CFG} is:

$S =$ "On input $\langle G, w \rangle$, where G is a CFG and w is a string:

1. Convert G to an equivalent grammar in Chomsky normal form
2. List all derivations with $2n - 1$ steps, $n = \text{length}(w)$ except if $n \leq 1$;
for $n \leq 1$ list all derivations with 1 step
3. If any of the derivations listed above generates w , *accept*, if not *reject*."

Observations

- The problem of testing whether G generates w is actually the parsing problem of compiling programming languages
- The algorithm performed by S is very inefficient. Early algorithm based on the same idea is $\mathcal{O}(n^3)$
- Theorem 2.20 proves that CFG are equivalent with PDA and provides a mechanism to convert a CFG into a PDA and vice-versa.

Conclusion: everything about decidability of problems concerning CFG applies equally to PDAs

Problem 2 Emptiness testing for CFGs

Problem: For a given CFG G is $L(G) = \emptyset$?

Language:

$$E_{CFG} = \{\langle G \rangle \mid G \text{ is a CFG and } L(G) = \emptyset\}$$

Theorem 4.8: E_{CFG} is a decidable language

Proof idea

- To test whether $L(G)$ is empty we need to test whether the axiom of G can generate a string of terminals
- We may solve however a more general problem, determining *for each variable* whether that variable can generate a string of terminals
- When the algorithm determines that a variable can generate a string of terminals the algorithm mark that variable
- The algorithm start by marking first all terminals. Then it marks variables that have on their rhs in some rules only terminals, i.e., marked symbols, and so one

Proof of theorem 4.8

Construct the TM R :

$R =$ "On input $\langle G \rangle$ where G is a CFG:

1. Mark all terminal symbols of G
2. Repeat until no new variable get marked:
Mark any variable A where G has a rule $A \rightarrow u_1 u_2 \dots u_k$ and each symbol u_1, u_2, \dots, u_k has already been marked
3. If the start symbol of G (i.e., the axiom) is not marked, *accept*; otherwise *reject*.

Problem 3 decidability of CFL

Problem: for an CFL A and string w does w belong to A ?, i.e. is there a CFG G such that $w \in L(G)$?

Language: $A_{CFL} = \{\langle A \rangle \mid A \text{ a CFL and } w \text{ string}\}$

Theorem 4.9 Every context-free language is decidable

Proof idea

Let A be a CFL

- A bad idea: convert a PDA for A directly into a TM.

Some branches of a PDA computation may go forever, reading and writing the stack without coming to a halt. The simulation TM would then have some non-halting branches in its computation and thus it would not be a decider

- A good idea: *use the TM S that decides string generation problem by converting G into Chomsky normal form*

Proof of Theorem 4.9

Let G be a CFG for A , i.e. $L(G) = A$. Design a TM M_G that decides A by building a copy of G into M_G :

$M_G =$ "On input w :

1. Run TM S on input $\langle G, w \rangle$
2. If this machine accepts, *accept*; if it rejects, *rejects*

Note: TM S converts G to Chomsky normal form, and produces all derivations of length $2n - 1$ where $n = |w|$. Then check if w is among the derived strings.

Fact 2

Class of CF languages is not closed under intersection

Proof: By construction.

- Consider the CF languages

$$A = \{a^m b^n c^n \mid m, n \geq 0\} \text{ and } B = \{a^n b^n c^m \mid m, n \geq 0\}$$

generated by the grammars:

$$S \rightarrow RT, R \rightarrow aR \mid \epsilon, T \rightarrow bTc \mid \epsilon \text{ and}$$

$$S \rightarrow TR, T \rightarrow aTb \mid \epsilon, R \rightarrow cR \mid \epsilon \text{ respectively.}$$

- $L(A) \cap L(B) = \{a^n b^n c^n \mid n \geq 0\}$ which is not a CFL
- This establishes Fact 2.

Fact 3

Class of CF languages is not closed under complementation

Proof: by contradiction.

Assume that CFL is closed under complementation

- If G_1 and G_2 are two CFG then $\overline{L(G_1)}$ and $\overline{L(G_2)}$ are CFL
- Then $\overline{L(G_1)} \cup \overline{L(G_2)}$ is a CFL. Hence, $\overline{\overline{L(G_1)} \cup \overline{L(G_2)}}$ is a CFL.
- By DeMorgan's law $\overline{\overline{L(G_1)} \cup \overline{L(G_2)}} = L(G_1) \cap L(G_2)$, a contradiction because class of CFL is not closed under intersections.

Problem 4 CFL equality problem

Problem: For two CFL languages generated by two CFGs

G and H is $L(G) = L(H)$ true?

Language: $EQ_{CFG} = \{\langle G, H \rangle \mid G, H \text{ CFGs and } L(G) = L(H)\}$

Note:

- Since class of CF languages is not closed under intersection and complementation (as seen before), we cannot use the symmetric difference for EQ_{CFG} .

- In fact EQ_{CFG} is not decidable (to be proven)

Methodology (review)

To solve decidability problems concerning relations between languages one should proceed as follows:

- Understand the relationship
- Transform the relationship into an expression using closure operators on decidable languages
- Design a TM that constructs the language thus expressed
- Run a TM that decide the language represented by the expression

Example 1

Equivalence of DFA and REX:

- Consider the problem of testing whether a DFA and a regular expression are equivalent.
- Express this problem as a language and show that this language is decidable

Solution

- Let $EQ_{DFA,REX} = \{\langle A, R \rangle \mid A \text{ is a DFA, } R \text{ is a regular expression and } L(A) = L(R)\}$
- The following TM E decides $EQ_{DFA,REX}$:
 $E =$ "On input $\langle A, R \rangle$
 1. Convert R to an equivalent DFA B
 2. Use the TM F for deciding EQ_{DFA} on input $\langle A, B \rangle$
 3. If F accepts, *accept*; if F rejects, *reject*."

Note: F constructs C , the DFA that recognizes the symmetric difference of A and B , $L(C) = L(A) \cap \overline{L(B)} \cup (\overline{L(A)} \cap L(B)$; and test if $L(C)$ is empty

Example 2

Decidability of Σ^*

- **Problem:** is the language Σ^* , for Σ a finite alphabet, decidable?
- **Language:** $ALL_{DFA} = \{\langle A \rangle \mid A \text{ is a DFA that recognize } \Sigma^*\}$
Show that ALL_{DFA} is decidable

Solution

The TM L that decides ALL_{DFA} uses the fact that $\overline{L(A)}$ is regular

L = "On input $\langle A \rangle$ where A is a DFA:

1. Construct DFA B that recognizes $\overline{L(A)}$ by swapping accept and unaccept states in A
2. Run the TM T that decides the emptiness E_{DFA} on B
3. If T accepts, *accept*; if T rejects *reject*."

Note: if T accepts it means that $L(B) = \emptyset$. But $L(B) = \overline{L(A)}$. That is,

$\Sigma^* \setminus L(A) = \emptyset$, i.e. $L(A) = \Sigma^*$.

Example 3

Using CFG and REX

- **Problem:** show that the problem of testing whether a CFG generates some string in 1^* is decidable.
- **Language:** $A = \{\langle G \rangle \mid G \text{ is a CFG over } \{0, 1\}^* \text{ and } 1^* \cap L(G) \neq \emptyset\}$

Solution

Assume that G is over $\{0, 1\}^*$. Then we need to show that the language $A = \{\langle G \rangle \mid G \text{ is a CFG over } \{0, 1\}^* \text{ and } 1^* \cap L(G) \neq \emptyset\}$ is decidable. Since 1^* is regular and $L(G)$ is CFL then $1^* \cap L(G)$ is a CFL. Hence the TM X that decides A is:

$X =$ "On input $\langle G \rangle$ where G is a CFG:

1. Construct CFG H such that $L(H) = 1^* \cap L(G)$
2. Run the TM R that decides the language E_{CFG} on $\langle H \rangle$
3. If R accepts, *reject*, if R rejects, *accept*."

Note: if R accepts it means that $L(H) = 1^* \cap L(G) = \emptyset$. That is, $\forall w \in 1^*$, $w \notin L(G)$, hence, X should reject.

Example 4

Example regular expressions

- **Problem:** Is the language generated by a particular regular expression decidable? For example, is the language of regular expressions that contain at least one string that has the pattern "111" as a substring decidable?
- **Language:** $A = \{\langle R \rangle \mid R \text{ is a regular expression describing a language that contain at least one string } w \text{ that has "111" as a substring (i.e., } w = x111y \text{ where } x \text{ and } y \text{ are strings)}\}$

Solution

The language A is decidable. The reason is that language A can be expressed using regular operators as $L(\Sigma^* \circ 111 \circ \Sigma^*) \cap L(R)$. Hence, the TM X that decides A is:

$X =$ "On input $\langle R \rangle$ where R is a regular expression:

1. Construct the DFA E that accepts $\Sigma^*111\Sigma^*$
2. Construct the DFA B that accepts $L(B) = L(R) \cap L(E)$
3. Run TM T that decide E_{DFA} on input $\langle B \rangle$
4. If T accepts *reject*, if T rejects *accept*."

Note: if T accepts, it means that $L(R) \cap L(E) = \emptyset$, i.e., $\forall x, y \in \Sigma^*$, $x111y \notin L(R)$.

The halting problem

- Our problem now is to test whether a Turing machine accepts a given input string.
- By analogy with A_{DFA} and A_{CFG} we call the corresponding language A_{TM}

$$A_{TM} = \{\langle M, w \rangle \mid M \text{ is a TM and } M \text{ accepts } w\}$$

- Contrasting A_{DFA} and A_{CFG} , which are decidable,
 A_{TM} is not decidable

Theorem 4.11

A_{TM} is undecidable

- A_{TM} is however Turing-recognizable
- Hence, Theorem 4.11 shows that recognizers are more powerful than deciders
- Requiring a TM to halt on all inputs restricts the kind of languages that it can recognize

A recognizer for A_{TM}

The following TM U recognizes A_{TM}

$U =$ "On input $\langle M, w \rangle$, where M is a TM and w is a string

1. Simulate M on input w
2. If M ever enters its accept state, *accept*, if M ever enters its reject state, *reject*".

Note

- Machine U loops on the input $\langle M, w \rangle$ if M loops on w . This is why U does not decide A_{TM} .
- If the algorithm had some way to determine that M was not halting on w , it could reject. This is why it is called the halting problem.
- However, as we demonstrate Theorem 4.11, an algorithm has no way to make this determination

Observations

1. The TM U is interesting in its own right because it is an example of the *universal Turing machine*, first proposed by Turing
2. U is called universal because it is capable to simulate any other Turing machine from the description of that machine
3. The universal TM played an important role in the theory of computation by stimulating the development of stored-program computers

Note: the algorithm implemented by a processor while executing a program:

```
while ((PC).opcode is not halt)
    Execute PC;
    PC := Next(PC);
```

behaves like U .