

1.3 Normal Forms of Context-Free Grammars

Up to now, we've imposed no restrictions on the form of each rule in our context-free grammars. As long as each rule of our context-free grammar looked like $A \rightarrow \alpha$, where A is a nonterminal symbol and α is a combination of terminal and nonterminal symbols, we were happy.

However, computers (and, by extension, the people who program computers) like having structure. For instance, a compiler for a programming language usually incorporates a context-free grammar into its workflow at some point during the compilation of a program, and having a highly-structured grammar makes the compiler's job of determining which rule to apply both easier and faster.

Therefore, in some cases, we might like to transform a context-free grammar into a *normal form*; that is, to modify the grammar in such a way that each rule follows a canonical form or template.

The normal form we will focus on, the *Chomsky normal form*, was (as the name suggests) first studied by Noam Chomsky in the late 1950s as he attempted to develop a model of natural language using grammars. A grammar in Chomsky normal form is such that each rule either has two nonterminal symbols or one terminal symbol on the right-hand side.

Definition 8 (Chomsky normal form). A context-free grammar is in Chomsky normal form if every rule in the grammar is of one of the two following forms:

1. $A \rightarrow BC$ for $A, B, C \in V$ with $B, C \neq S$; or
2. $A \rightarrow a$ for $A \in V$ and $a \in \Sigma$.

Additionally, we may allow the rule $S \rightarrow \epsilon$.

The main benefit of converting a grammar into Chomsky normal form comes in how we can represent and store derivations of words in memory. Since each rule derives either two nonterminal symbols or one terminal symbol, every parse tree will have a branching factor of either 2 or 1. This fact allows us to use efficient data structures for representing binary trees in memory, as well as to apply efficient algorithms to process parse trees and derivations. Moreover, the number of steps in a derivation using a grammar in Chomsky normal form is easy to bound: if the grammar generates a word w , then the derivation of w will contain $|w| - 1$ applications of a rule of the first form and $|w|$ applications of a rule of the second form.

Example 9. Let $\Sigma = \{a, b\}$, and consider the following two grammars. Each grammar generates words consisting of one b surrounded on either side by zero or more a s. The grammar on the left is not in Chomsky normal form. The grammar on the right is in Chomsky normal form, and it is equivalent to the grammar on the left.

$$\begin{array}{ll}
 S \rightarrow AbA & S_0 \rightarrow TA \mid BA \mid AB \mid b \\
 A \rightarrow Aa \mid \epsilon & T \rightarrow AB \\
 & A \rightarrow AC \mid a \\
 & B \rightarrow b \\
 & C \rightarrow a
 \end{array}$$

Every context-free grammar can be converted into a context-free grammar in Chomsky normal form, and the conversion process consists of four steps:

1. **START:** Replace the start nonterminal.

Add a new start nonterminal S_0 together with a new rule $S_0 \rightarrow S$, where S is the start nonterminal of the original grammar. This guarantees that S_0 will not occur on the right-hand side of any rule.

2. **EPS:** Remove epsilon rules.

Remove all rules of the form $A \rightarrow \epsilon$, where $A \neq S$. For each occurrence of A on the right-hand side of a rule in the original grammar, add a new rule with that occurrence of A removed from the right-hand side; that is, convert all rules of the form $X \rightarrow \alpha A \beta$ to $X \rightarrow \alpha \beta$, where α and β are combinations

of nonterminal and terminal symbols. Note that this applies to each *occurrence* of A , so a rule of the form $X \rightarrow \alpha A \beta A \gamma$ would contribute three new rules: $X \rightarrow \alpha \beta A \gamma$, $X \rightarrow \alpha A \beta \gamma$, and $X \rightarrow \alpha \beta \gamma$.

If there exists a rule of the form $X \rightarrow A$, then add a new rule of the form $X \rightarrow \epsilon$ unless we previously removed a rule of that form.

Repeat until all epsilon rules not involving the original start nonterminal are removed.

3. **UNIT:** Remove unit rules.

Unit rules are rules of the form $A \rightarrow B$, where A and B are nonterminal symbols. Remove all rules of this form.

If there exists a rule of the form $B \rightarrow \alpha$, where α is a combination of nonterminal and terminal symbols, then add a new rule $A \rightarrow \alpha$ unless we previously removed a unit rule of that form.

Repeat until all unit rules are removed.

4. **BIN:** Remove rules with more than two nonterminal or terminal symbols on the right-hand side.

Replace each rule of the form $A \rightarrow \alpha_1 \alpha_2 \dots \alpha_k$, where $k \geq 3$ and each α_i is either a nonterminal or terminal symbol, with a series of new rules $A \rightarrow \alpha_1 A_1$, $A_1 \rightarrow \alpha_2 A_2$, \dots , $A_{k-2} \rightarrow \alpha_{k-1} \alpha_k$. Each A_i is a new nonterminal symbol.

If α_i is a terminal symbol in one of our new rules $A_{i-1} \rightarrow \alpha_i A_i$, then remove this rule, replace α_i with a new nonterminal symbol B_i , and add two new rules $B_i \rightarrow \alpha_i$ and $A_{i-1} \rightarrow B_i A_i$.

Example 10. Consider the following grammar not in Chomsky normal form:

$$\begin{aligned} S &\rightarrow ASB \\ A &\rightarrow \mathbf{a}AS \mid \mathbf{a} \mid \epsilon \\ B &\rightarrow S\mathbf{b}S \mid A \mid \mathbf{bb} \end{aligned}$$

We will convert this grammar to an equivalent grammar in Chomsky normal form.

1. **START:** Replace the start nonterminal.

Adding a new start nonterminal and the rule $S_0 \rightarrow S$ gives us the following:

$$\begin{aligned} S_0 &\rightarrow S \\ S &\rightarrow ASB \\ A &\rightarrow \mathbf{a}AS \mid \mathbf{a} \mid \epsilon \\ B &\rightarrow S\mathbf{b}S \mid A \mid \mathbf{bb} \end{aligned}$$

2. **EPS:** Remove epsilon rules.

The first epsilon rule we will remove is $A \rightarrow \epsilon$. Since A occurs on the right-hand side of two other rules ($S \rightarrow ASB$ and $A \rightarrow \mathbf{a}AS$), removing this rule means we must add two new rules $S \rightarrow SB$ and $A \rightarrow \mathbf{a}S$. Additionally, since we have a rule $B \rightarrow A$, we must add a new rule $B \rightarrow \epsilon$. This gives the following:

$$\begin{aligned} S_0 &\rightarrow S \\ S &\rightarrow ASB \mid SB \\ A &\rightarrow \mathbf{a}AS \mid \mathbf{a} \mid \mathbf{a}S \\ B &\rightarrow S\mathbf{b}S \mid A \mid \mathbf{bb} \mid \epsilon \end{aligned}$$

Next, we remove the epsilon rule $B \rightarrow \epsilon$ that we just added. Since B occurs on the right-hand side of two other rules ($S \rightarrow ASB$ and $S \rightarrow SB$), removing this rule means we must add two new rules $S \rightarrow AS$ and $S \rightarrow S$. This gives the following:

$$\begin{aligned}
 S_0 &\rightarrow S \\
 S &\rightarrow ASB \mid SB \mid AS \mid S \\
 A &\rightarrow aAS \mid a \mid aS \\
 B &\rightarrow SbS \mid A \mid bb
 \end{aligned}$$

3. **UNIT:** Remove unit rules.

We first remove the unit rule $B \rightarrow A$. We do so by replacing A on the right-hand side of the rule with everything that can be produced by a rule of the form $A \rightarrow \alpha$:

$$\begin{aligned}
 S_0 &\rightarrow S \\
 S &\rightarrow ASB \mid SB \mid AS \mid S \\
 A &\rightarrow aAS \mid a \mid aS \\
 B &\rightarrow SbS \mid bb \mid aAS \mid a \mid aS
 \end{aligned}$$

Next, we remove the unit rule $S \rightarrow S$. This change is more straightforward, as we don't need to modify the S rule in any other way:

$$\begin{aligned}
 S_0 &\rightarrow S \\
 S &\rightarrow ASB \mid SB \mid AS \\
 A &\rightarrow aAS \mid a \mid aS \\
 B &\rightarrow SbS \mid bb \mid aAS \mid a \mid aS
 \end{aligned}$$

Lastly, we remove the unit rule $S_0 \rightarrow S$. Again, we replace S on the right-hand side of the rule with everything that can be produced by a rule of the form $S \rightarrow \alpha$:

$$\begin{aligned}
 S_0 &\rightarrow ASB \mid SB \mid AS \\
 S &\rightarrow ASB \mid SB \mid AS \\
 A &\rightarrow aAS \mid a \mid aS \\
 B &\rightarrow SbS \mid bb \mid aAS \mid a \mid aS
 \end{aligned}$$

4. **BIN:** Remove rules with more than two nonterminal or terminal symbols on the right-hand side.

Starting from the top, we replace the rule $S_0 \rightarrow ASB$ with the rules $S_0 \rightarrow AU_1$ and $U_1 \rightarrow SB$:

$$\begin{aligned}
 S_0 &\rightarrow AU_1 \mid SB \mid AS \\
 S &\rightarrow ASB \mid SB \mid AS \\
 A &\rightarrow aAS \mid a \mid aS \\
 B &\rightarrow SbS \mid bb \mid aAS \mid a \mid aS \\
 U_1 &\rightarrow SB
 \end{aligned}$$

Similarly, we replace the rule $S \rightarrow ASB$ with the rules $S \rightarrow AU_2$ and $U_2 \rightarrow SB$:

$$\begin{aligned}
 S_0 &\rightarrow AU_1 \mid SB \mid AS \\
 S &\rightarrow AU_2 \mid SB \mid AS \\
 A &\rightarrow aAS \mid a \mid aS \\
 B &\rightarrow SbS \mid bb \mid aAS \mid a \mid aS \\
 U_1 &\rightarrow SB \\
 U_2 &\rightarrow SB
 \end{aligned}$$

Next, we replace the rule $A \rightarrow aAS$ with the rules $A \rightarrow aU_3$ and $U_3 \rightarrow AS$:

$$\begin{aligned}
 S_0 &\rightarrow AU_1 \mid SB \mid AS \\
 S &\rightarrow AU_2 \mid SB \mid AS \\
 A &\rightarrow aU_3 \mid a \mid aS \\
 B &\rightarrow SbS \mid bb \mid aAS \mid a \mid aS \\
 U_1 &\rightarrow SB \\
 U_2 &\rightarrow SB \\
 U_3 &\rightarrow AS
 \end{aligned}$$

Moving along, we replace the rules $B \rightarrow SbS$ and $B \rightarrow aAS$ with the set of rules $B \rightarrow SU_4$, $B \rightarrow aU_5$, $U_4 \rightarrow bS$, and $U_5 \rightarrow AS$:

$$\begin{aligned}
 S_0 &\rightarrow AU_1 \mid SB \mid AS \\
 S &\rightarrow AU_2 \mid SB \mid AS \\
 A &\rightarrow aU_3 \mid a \mid aS \\
 B &\rightarrow SU_4 \mid bb \mid aU_5 \mid a \mid aS \\
 U_1 &\rightarrow SB \\
 U_2 &\rightarrow SB \\
 U_3 &\rightarrow AS \\
 U_4 &\rightarrow bS \\
 U_5 &\rightarrow AS
 \end{aligned}$$

Lastly, we must replace all rules containing one terminal and one nonterminal symbol on the right-hand side. We introduce two new rules $V_1 \rightarrow a$ and $V_2 \rightarrow b$ and make the appropriate changes to other rules:

$$\begin{aligned}
 S_0 &\rightarrow AU_1 \mid SB \mid AS \\
 S &\rightarrow AU_2 \mid SB \mid AS \\
 A &\rightarrow V_1U_3 \mid a \mid V_1S \\
 B &\rightarrow SU_4 \mid V_2V_2 \mid V_1U_5 \mid a \mid V_1S \\
 U_1 &\rightarrow SB \\
 U_2 &\rightarrow SB \\
 U_3 &\rightarrow AS \\
 U_4 &\rightarrow V_2S & V_1 &\rightarrow a \\
 U_5 &\rightarrow AS & V_2 &\rightarrow b
 \end{aligned}$$

1.4 Context-Free Languages

With our notion of a context-free grammar, it's easy for us to define a context-free language. Just like we defined regular languages in terms of regular operations or regular expressions, we have the following definition for context-free languages in terms of context-free grammars.

Definition 11 (Context-free language). If some language L is generated by a context-free grammar, then L is context-free.

Thus, each of the languages in this lecture for which we constructed context-free grammars—the language of words $a^n b^n$, the language of balanced parentheses, the language of words $a^i b^j c^k$ where $i = j$ or $j = k$ —are context-free languages. Our definition also tells us that a language is context-free if we can construct a context-free grammar generating words in that language. However, be careful: any such context-free grammar must generate *all* and *only all* the words in the given language.

As a shorthand, we denote the class of languages generated by a context-free grammar by CFG.

Example 12. Consider the language $L = \{a^{2i} b^i c^{j+2} \mid i, j \geq 0\}$ over the alphabet $\Sigma = \{a, b, c\}$. Here, we can see that the counts of *as* and *bs* are related, while the number of *cs* is independent of the number of *as* and *bs*.

Let's construct a context-free grammar generating words in L . Evidently, we will need two kinds of rules: one rule will generate the *as* and *bs* together, while the other rule will generate the *cs*. We can use the start nonterminal to apply these rules in the correct order. Our context-free grammar will therefore look like the following:

$$\begin{aligned} S &\rightarrow UV \\ U &\rightarrow \mathbf{aaUb} \mid \epsilon \\ V &\rightarrow \mathbf{cV} \mid \mathbf{cc} \end{aligned}$$

Let's now take a look at each rule in turn. The first rule, S , ensures that we apply the U rule before the V rule. This in turn ensures that all *as* and *bs* occur before the *cs* in the generated word. The second rule, U , either recursively produces two *as* and one *b* or produces the empty word. This ensures that we maintain the correct count of $2i$ *as* and i *bs*. Finally, the third rule, V , either recursively produces one *c* or produces the symbols *cc*. This ensures that we have exactly $j + 2$ *cs* in our generated word.

Example 13. Recall our context-free language $L_{a=b} = \{a^n b^n \mid n \geq 1\}$. Here, let's consider a more general language: $L_{\text{mixed}a=b} = \{w \in \{a, b\}^* \mid |w|_a = |w|_b\}$. Observe that the main difference with this language is that the order of *as* and *bs* no longer matters; we just need the same count of *as* and *bs*. Can we construct a context-free grammar for $L_{\text{mixed}a=b}$?

Since order no longer matters, we just need our context-free grammar to generate a pair of *as* and *bs* each time we add terminal symbols. For this, we can use essentially the same rule as we used in our context-free grammar for $L_{a=b}$: $S \rightarrow \mathbf{aSb} \mid \mathbf{bSa}$. We also need a rule that allows us to mix the order of *as* and *bs*; for instance, to place two *as* or two *bs* next to each other instead of strictly nesting pairs. For this, we can use a rule similar to one we included in our context-free grammar for L_{\cup} : $S \rightarrow \mathbf{SS}$.

Thus, our context-free grammar will look like the following:

$$S \rightarrow \mathbf{SS} \mid \mathbf{aSb} \mid \mathbf{bSa} \mid \epsilon$$

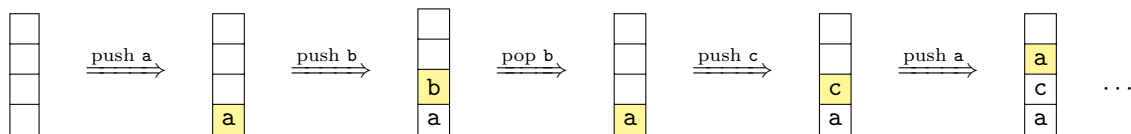
2 Pushdown Automata

When we first introduced finite automata as a computational model for regular languages, we emphasized the facts that finite automata have no memory, no storage, and no ability to return to a previously-read symbol. Naturally, these restrictions limited the kinds of languages the model is able to recognize, and we showed that such restrictions resulted in the model recognizing exactly the class of regular languages.

At the end of the previous lecture, we saw that there exist languages that are not regular, and therefore are not recognized by finite automata. We know now that the next “step” of our language hierarchy is the class of context-free languages. Thus, a new question arises: what kind of computational model is capable of recognizing context-free languages?

Since every context-free language is generated by a context-free grammar, and since we know that context-free grammars must “remember” which nonterminal and terminal symbols are being manipulated over the course of a derivation, any model of computation recognizing context-free languages must include a form of memory. What is the best form of memory to use in this situation? If we view a derivation as a parse tree, then the derivation progresses as we go deeper into the parse tree, and we can easily model the depth of a derivation using *stack* memory.

As a brief review, a stack is a data structure with two operations that manipulate data: *push* and *pop*. Pushing a symbol to a stack adds it to the top of the stack and moves all other symbols one position down in the stack. Conversely, popping a symbol from a stack removes it from the top and moves all other symbols one position up. As a result, a stack provides last-in-first-out (LIFO) storage.⁴ We can view the symbol at the top of the stack at any time during a computation, but we cannot view any other symbols in the stack unless we pop the symbol currently at the top of the stack.



Since we’re dealing with an abstract model of computation and not a real-world computer, we can make the assumption that our stack size is *unbounded*; that is, we can push as many symbols to the stack as we want without worrying about running out of space.

Now that we have our form of storage established, we can define our model of computation. At its core, this model is a finite automaton with a stack added to it. Since the automaton is now able to push symbols to a stack, we give it an appropriate name: a *pushdown automaton*.⁵

In addition to reading a symbol of its input word on a transition, a pushdown automaton can read from and write to the stack on the same transition. In order to perform this double duty, we specify two alphabets for a pushdown automaton: the *input alphabet*, which contains symbols used in the input word, and the *stack alphabet*, which contains symbols the pushdown automaton can use in its stack. This allows us to combine actions on the input word and actions on the stack in a single transition, without risking confusion over the meaning of any particular alphabet symbol. The transitions of a pushdown automaton may additionally use epsilon for either the input word action (i.e., not reading a symbol of the input word) or the stack action (i.e., not pushing to/popping from the stack). Just like we denoted a finite automaton’s alphabet by Σ , we will use Σ to denote a pushdown automaton’s input alphabet. Likewise, we will use Γ to denote the stack alphabet.

In order for our model of computation to use two alphabets at once, we must modify its transition function accordingly. Recall that a finite automaton (with epsilon transitions) transitions on a pair (q, a) , where $q \in Q$ and $a \in \Sigma \cup \{\epsilon\}$. By comparison, a pushdown automaton transitions on a tuple (q, a, b) , where $q \in Q$, $a \in \Sigma \cup \{\epsilon\}$, and $b \in \Gamma \cup \{\epsilon\}$. Thus, a pushdown automaton uses both the current symbol of its input word (or epsilon) as well as the top symbol of its stack (or epsilon) to determine to which state it will transition. After transitioning, the pushdown automaton will be in a possibly-different state, and it will have a possibly-different symbol at the top of its stack.

Lastly, a pushdown automaton has no mechanism for detecting whether or not its stack is empty. To avoid any potential issues, we often incorporate a special “bottom of stack” symbol \perp into the transitions of a

⁴By comparison, a data structure like a queue would provide first-in-first-out (FIFO) storage.

⁵The name “pushdown automaton” doesn’t specifically come from its ability to push symbols, but rather from an older term for a stack: a *pushdown store*.