

Lecture 14: Repetition in context free languages

10 March 2009

1 Generating new words

We are interested in phenomena of repetition in context free languages. We had seen that regular languages repeat themselves if the strings are sufficiently long. We would like to make a similar statement about regular languages, but unfortunately, while the general statement is correct, the details are somewhat more involved.

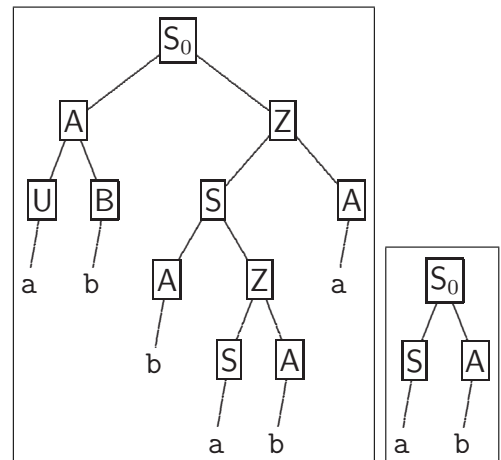
1.1 Example of repetition

As a concrete example, consider the following context-free grammar which is in Chomsky normal form (CNF). (We remind the reader that any context free grammar can be converted into CNF, as such assuming that we have a grammar in CNF form does not restrict our discussion.)

As a concrete example, consider the following grammar from the previous lecture:

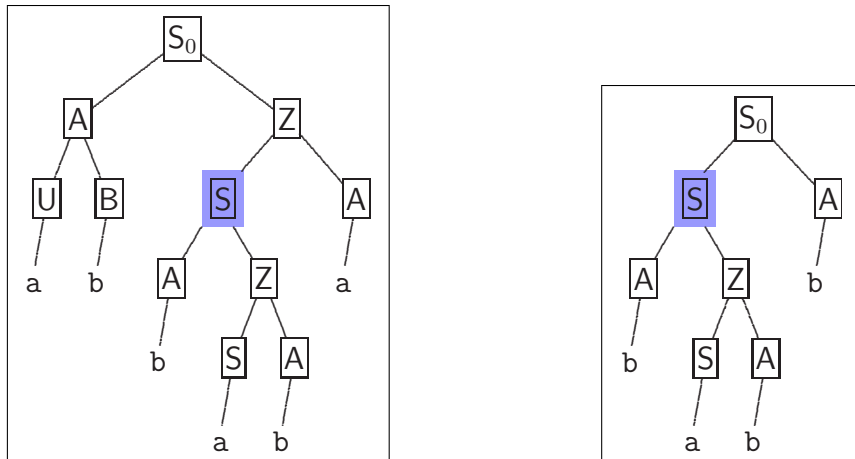
$$\begin{aligned}
 \Rightarrow \quad & S_0 \rightarrow AZ \mid UB \mid a \mid SA \mid AS \\
 & S \rightarrow AZ \mid UB \mid a \mid SA \mid AS \\
 & A \rightarrow b \mid AZ \mid UB \mid a \mid SA \mid AS \\
 & B \rightarrow b \\
 & U \rightarrow a \\
 & Z \rightarrow SA
 \end{aligned}$$

(G5)

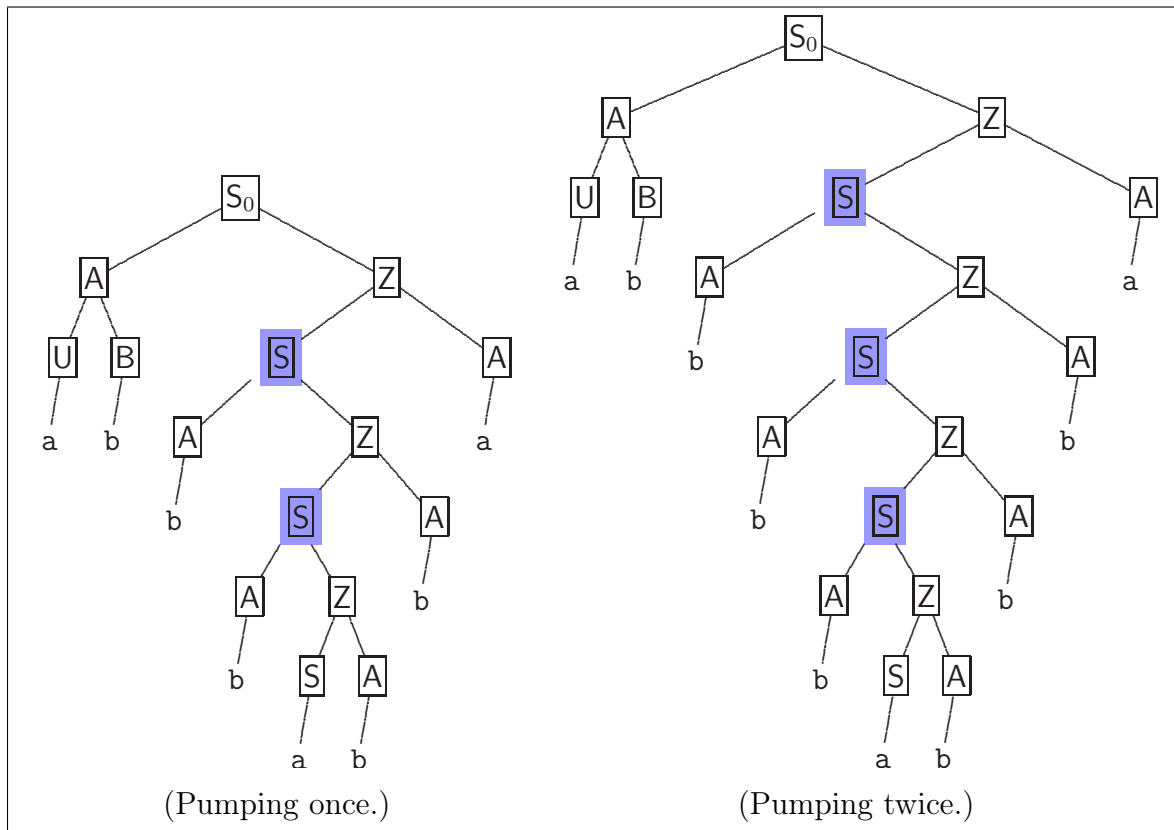


Next, consider the two words created from this grammar, as depicted on the right.

What if we wanted to make the second word longer without thinking too much? Well, both parse trees have subtrees with nodes generated by the variable S . As such, we could just cut the subtree of the first word using the variable S , and replace the subtree of S in the second word by this subtree, which would look like the following:



Even more interestingly, we can do this cut and paste on the original tree:



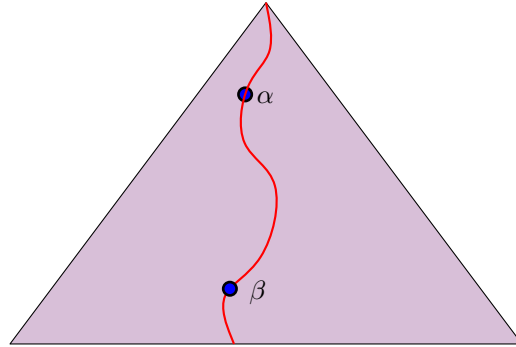
Naturally, we can repeat this pumping operation (cutting and pasting a subtree) as many times as want, see for example Figure 1. In particular, we get that the word $abb^i ab^i a$, for any i , is in the language of the grammar (G_5) . Notice that unlike the pumping lemma for regular languages, here the repetition happens in two places in the string. We claim that such a repetition (in two places) in the word must happen for any context free language, once we take a word which is sufficiently long.

2 The pumping lemma for CFG languages

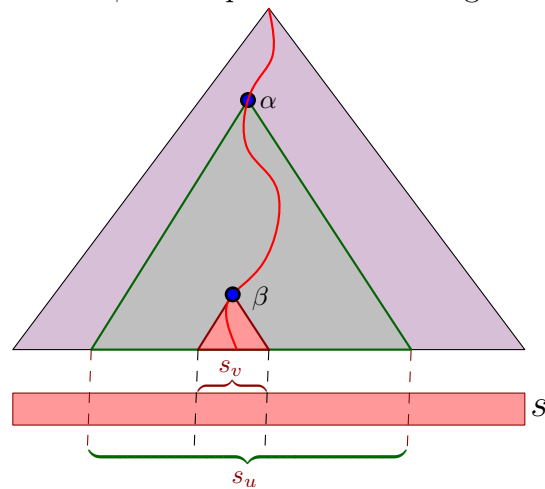
So, assume we are given a context free grammar \mathcal{G} which is in CNF, and it has m variables.

2.1 If a variable repeats

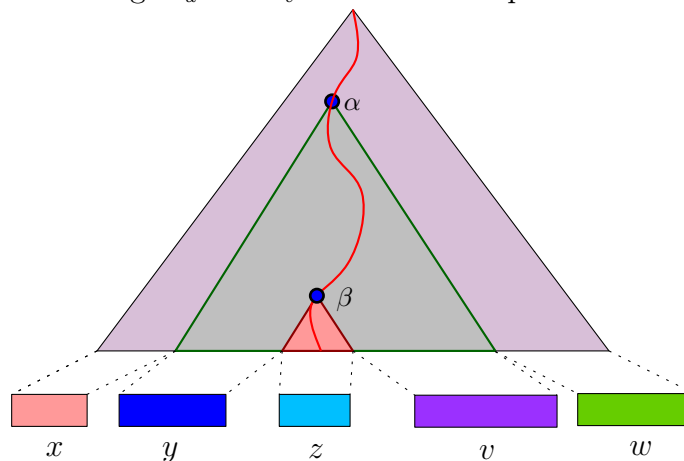
So, assume we have a parsing tree T for a word s (where the underlying grammar is in CNF), and there is a path in T from the root to a leaf, such that a variables repeats twice. So, say nodes α and β have the same variable (say S) stored in them:



The subtrees rooted at α and β corresponds to substrings of s :

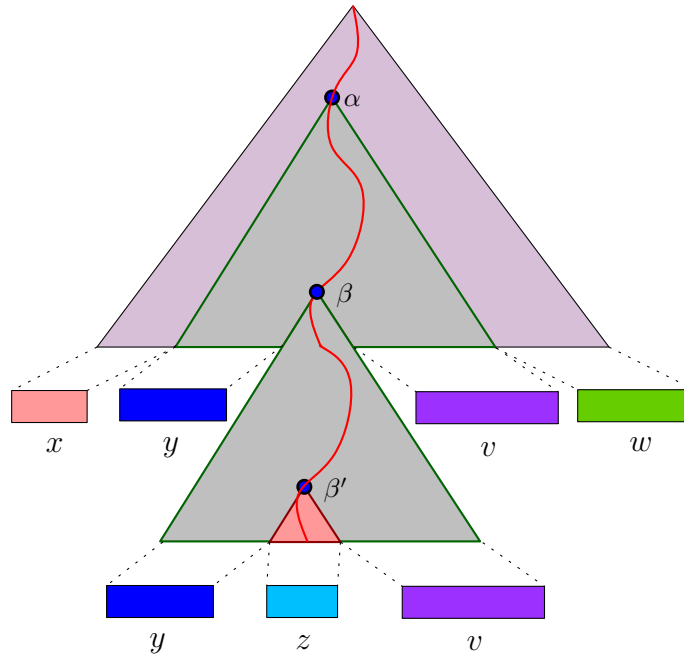


In particular, the substrings s_u and s_v break s into 5 parts:



Namely, s can be written as $s = xyzvw$.

Now, if we copy the subtree rooted at α and copy it to β , we get a new parse tree:

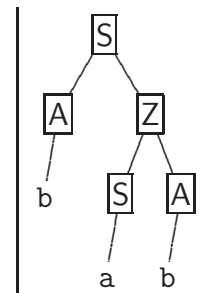


The new tree is much bigger, and the new string it represents is $s = xyzvzw$. In general, if we do this cut & paste operation $i - 1$ times, we get the string

$$xy^i z v^i w.$$

2.2 How tall the parse tree have to be?

We will refer to a parse generated from a context free grammar in CNF form as a **CNF tree**. A CNF tree has the special property that the parent of a leaf as a single child which is a terminal. The **height** of a tree is the maximum number of edges on a path from the root of the tree to a leaf. Thus, the tree depicted on the right has height 3.



The grammar \mathcal{G} has m variables. As such, if the parse tree T has a path π from the root of length k , and $k > m$ (i.e., the path has k edges), then it must contain at least $m + 1$ variables (the last edge is between a variable and a terminal). As such, by the pigeon hole principle, there must be a repeated variable along π . In particular, a parse tree that does not have a repeated variable have height at most m .

Since \mathcal{G} is in CNF, its a binary tree, and a variable either has two children, or a single child which is a leaf (and that leaf contains a single character of the input). As such, a tree of height at most m , contains at most 2^m leaves¹, and represents as such a string of length at most 2^m .

¹In fact, a CNF tree of height m can have at most 2^{m-1} leaves (figure out why), but thats a subtlety we will ignore that anyway works in our favor.

We restate the above observation formally for the record.

Observation 2.1 *If a CNF parse tree (or a subtree of such a tree) has height h , then the string it generates is of length at most 2^h .*

Lemma 2.2 *Let \mathcal{G} be a grammar given in Chomsky Normal Form (CNF), and consider a word $s \in \mathbf{L}(\mathcal{G})$, such that $\ell = |s|$ is strictly larger than 2^m (i.e., $\ell > 2^m$). Then, any parse tree T for s (generated by \mathcal{G}) must have a path from the root to some leaf with a repeated variable on it.*

Proof: Assume for the sake of contradiction that T has no repeated variable on any path from the root, then the height of T is at most m . But a parse tree of height m for a CNF can generate a string of length at most 2^m . A contradiction, since $\ell = |s| > 2^m$. ■

2.3 Pumping Lemma for CNF grammars

We need the following observation.

Lemma 2.3 (CNF is effective.) *In a CNF parse tree T , if u and v are two nodes, both storing variables in them, and u is an ancestor of v , then the string S_u generated by the subtree of u is strictly longer than the substring S_v generated by the subtree of u . Namely, $|S_u| > |S_v|$. (Of course, S_v is a substring of S_u .)*

Proof: Assume that the node u stores a variable X , and that we had used the rule $X \rightarrow BC$ to generate its two children u_L and u_R . Furthermore, assume that u_L and u_R generated the strings S_L and S_R , respectively. The string generated by v must be a substring of either S_L or S_R . However, CNF has the property that no variable² can generate the empty word ϵ . As such $|S_R| > 0$ and $|S_L| > 0$.

In particular, assume without loss of generality, that v is in the left subtree of u , and as such S_v is a substring of S_L . We have that

$$|S_v| \leq |S_L| < |S_L| + |S_R| = |S_u|.$$

Lemma 2.4 *Let T be a tree, and π be the longest path in the tree realizing the height h of T . Fix $k \geq 0$, and let u be the k th node from the end of π (i.e., u is in distance $h - k$ from the root of T). Then the tree rooted at u has height at most k .*

Proof: Let r be the root of T , and assume, for the sake of contradiction, that T_u (i.e., the subtree rooted at u) has height larger than k , and let σ be the path from u to the leaf γ of T_u realizing this height (i.e., the length of σ is $> k$). Next, consider the path formed by concatenating the path in T from r to u with the path σ . Clearly, this is a new path of length $h - k + |\sigma| > h$ that leads from the root of T into a leaf of T . As such, the height of T is larger than h , which is a contradiction. ■

²Except the start variable, but this not relevant here.

Lemma 2.5 (Pumping lemma for Chomsky Normal Form (CNF).) *Let \mathcal{G} be a CNF context-free grammar with m variables in it. Then, given any word \mathcal{S} in $L(\mathcal{G})$ of length $> 2^m$, one can break \mathcal{S} into 5 substrings $\mathcal{S} = xyzvw$, such that for any $i \geq 0$, we have that xy^izv^iw is a word in $L(\mathcal{G})$. In addition, the following holds:*

1. *The strings y and v are not both empty (i.e., the pumping is getting us new words).*
2. *$|yzv| \leq 2^m$.*

Proof: Let T be a CNF parse tree for \mathcal{S} (generated by \mathcal{G}). Since $\ell = |\mathcal{S}| > 2^m$, by Lemma 2.2, there is a path in T from its root to a leaf which has a repeated variable (and its length is longer than m). In fact, let π be the longest path in T from the root to a leaf (i.e., π is the path realizing the height of the tree T). We know that T has more than $m + 1$ variables on it and as such it has a repetition.

We need to be a bit careful in picking the two nodes α and β on π to apply the pumping to. In particular, let α be the last node on π such that there is a repeated appearance of the symbol stored in u later in the path. Clearly, the length of the subpath τ of π starting at α till the end of π has at most m symbols on it (because otherwise, there would be another repetition on π). Let β be the node of $\tau \subseteq \pi$ which has repetition of the symbol stored in α .

By Lemma 2.4 the subtree T_α (i.e., the subtree of T rooted at α) has height at most m . As above, T_α and T_β generate two strings \mathcal{S}_α and \mathcal{S}_β , respectively. By Observation 2.1, we have that $|\mathcal{S}_\alpha| \leq 2^m$. By Lemma 2.3, we have that $|\mathcal{S}_\alpha| > |\mathcal{S}_\beta|$. As such, the two substrings \mathcal{S}_α and \mathcal{S}_β breaks \mathcal{S} into 5 substrings $\mathcal{S} = xyzvw$. Here, we have

$$\mathcal{S} = x \underbrace{y \quad z \quad v}_{=\mathcal{S}_\beta}^{\mathcal{S}_\alpha} w.$$

As such, we know that $|yv| = |\mathcal{S}_\alpha| - |\mathcal{S}_\beta| > 0$. Namely, the strings y and v are not both empty. Furthermore, $|yzv| = |\mathcal{S}_\alpha| \leq 2^m$.

The remaining task is to show the pumping. Indeed, if we replace T_β by the tree T_α we get a parse tree generating the string xy^2zv^2w . If we repeat this process $i - 1$ times, we get the word

$$xy^izv^iw \in L(\mathcal{G}),$$

for any i , establishing the lemma. ■

Lemma 2.6 (Pumping lemma for context-free languages.) *If L is a context-free language, then there is a number p (the pumping length) where, if \mathcal{S} is any string in L of length at least p , then \mathcal{S} may be divided into five pieces $\mathcal{S} = xyzvw$ satisfying the conditions:*

1. *for any $i \geq 0$, we have $xy^izv^iw \in L$,*
2. *$|yv| > 0$,*
3. *and $|yzv| \leq p$.*

Proof: Since L is context free it has a CNF grammar \mathcal{G} that generates it. Now, if m is the number of variables in \mathcal{G} , then for $p = 2^m + 1$, the lemma follows by Lemma 2.5. ■

3 Languages which are not context-free

3.1 The language $a^n b^n c^n$ is not context-free

Lemma 3.1 *The language $L = \{a^n b^n c^n \mid n \geq 0\}$ is not context-free.*

Proof: Assume, for the sake of contradiction, that L is context-free, and apply the Pumping Lemma to it (Lemma 2.6). As such, there exists $p > 0$ such that any word in L longer than p can be pumped. So, consider the word $\mathcal{S} = a^{p+1} b^{p+1} c^{p+1}$. By the pumping lemma, it can be written as $a^{p+1} b^{p+1} c^{p+1} = xyzvw$, where $|yzv| \leq p$.

We claim, that yzv can be made out of only two characters. Indeed, if yzv contained all three characters, it would have to contain the string b^{p+1} as a substring (as b^{p+1} separates all the appearances of a from all the appearances of c in \mathcal{S}). This would require that $|yzv| > p$ but we know that $|yzv| \leq p$.

In particular, let i_a, i_b and i_c be the number of a 's, b 's and c 's in the string yzv , respectively. All we know is that $i_a + i_b + i_c = |yzv| > 0$ and that $i_a = 0$ or $i_c = 0$. Namely, $i_a \neq i_b$ or $i_b \neq i_c$ (the case $i_a \neq i_c$ implies one of these two cases). In particular, by the pumping lemma, the word

$$\mathcal{S}_2 = xy^2zv^2w \in L.$$

We have the following:

character	how many times it appears in \mathcal{S}_2
a	$p + 1 + i_a$
b	$p + 1 + i_b$
c	$p + 1 + i_c$

If $i_a \neq i_b$ then \mathcal{S}_2 , by the above table, does not have the same number of a 's and b 's and as such it is not in L .

If $i_b \neq i_c$ then \mathcal{S}_2 , by the above table, does not have the same number of b 's and c 's and as such it is not in L .

In either case, we get that $\mathcal{S}_2 \notin L$, which is a contradiction. Namely, our assumption that L is context-free is false. ■

4 Closure properties

4.1 Context-free languages are not closed under intersection

We know that the languages

$$L_1 = \{a^* b^n c^n \mid n \geq 0\} \text{ and } L_2 = \{a^n b^n c^* \mid n \geq 0\}$$

are context-free (prove this). But

$$L = \{a^n b^n c^n \mid n \geq 0\} = L_1 \cap L_2$$

is not context-free by Lemma 3.1. We conclude that the intersection of two context-free languages is not necessarily context-free.

Lemma 4.1 *Context-free languages are not closed under intersection.*

4.2 Context-free languages are not closed under complement

Lemma 4.2 *Context-free languages are not closed under complement.*

Proof: Consider the language

$$L_3 = \left\{ a^i b^j c^k \mid i \neq j \text{ or } j \neq k \right\}.$$

The language L_3 is clearly context-free (why?). Consider its complement language $\overline{L_3}$. If L_3 is context-free, and context-free languages are closed under complement (which we are assuming here for the sake of contradiction), then $\overline{L_3}$ is context-free.

Now, $\overline{L_3}$ contains many strings we are not interested in (for example, $ccccba$). So, let us intersect it with the regular language $a^*b^*c^*$. We proved (in previous lecture), that the intersection of a context-free language and a regular language is still context-free. As such,

$$\widehat{L} = \overline{L_3} \cap \{a^*b^*c^*\}$$

is context-free. But this language contains all the strings $a^i b^j c^k$, where $i = j$ and $j = k$. That is, it's the language of Lemma 3.1, which we know is not context-free. A contradiction. ■

Proof: (Alternative proof.) The language $L = a^n b^n c^n$ can be written, by De Morgan laws, as

$$L = a^n b^n c^* \cap a^* b^n c^n = \overline{\overline{a^n b^n c^*} \cup \overline{a^* b^n c^n}}. \quad (1)$$

We assume for the sake of contradiction, that context-free languages are closed under complement, and we already know that they are closed under union. However, the languages

$$a^n b^n c^* \text{ and } a^* b^n c^n$$

are context-free. As such, by closure properties and Eq. (1), the language L is context-free, which is a contradiction to Lemma 3.1. ■