

# Discussion : Context-Free Grammars

March 3, 2009

## Questions on homework or exam?

Any questions? Complaints, etc?

## 1 Context free grammar for languages with balance or without it

### 1.1 Balanced strings

Let  $L_{a=b} = \{a^n b^n \mid n \geq 1\}$ . Here is a grammar for this language

$$S \rightarrow aSb \mid \epsilon.$$

### 1.2 Mixed balanced strings

Let  $L_{a=b}^{\text{mix}}$  be the language of all strings over  $\{a, b\}$ , with equal number of as and bs, where the as and bs might be mixed together.

We the following grammar for generating all strings that have the same numbers of  $a$ 's and  $b$ 's

$$S \rightarrow aSb \mid bSa \mid SS \mid \epsilon,$$

where  $S$  is the start variable.

**Advice To TA::** State the following lemma, and sketch its proof, but do not do the proof in the discussion section. Point the interested students to the class notes. **end**

**Lemma 1.1** We have  $L(\mathcal{G}) = L_{a=b}^{\text{mix}}$ .

*Proof:* It is easy to see that every string that  $\mathcal{G}$  generates has equal number of  $a$ 's and  $b$ 's. As such,  $L(\mathcal{G}) \subseteq L_{a=b}^{\text{mix}}$ .

We will use induction on the length of string  $x \in L(\mathcal{G})$ ,  $2n = |x|$ . For  $n = 0$  we can generate  $\epsilon$  by  $\mathcal{G}$ . For  $n = 1$ , we can generate both  $ab$  and  $ba$  by  $\mathcal{G}$ .

Now for  $n > 1$ , consider a balanced string with length  $2n$ ,  $x = x_1 x_2 x_3 \cdots x_{2n} \in L_{a=b}^{\text{mix}}$ . Let  $\#_c(y)$  be the number of appearances of the character  $c$  in the string  $y$ . Let  $\alpha_i = \#_a(x_1 \cdots x_i) - \#_b(x_1 \cdots x_i)$ . Observe that  $\alpha_0 = \alpha_{2n} = 0$ . If  $\alpha_j = 0$ , for some  $1 < j < 2n$ ,

then we can break  $x$  into two words  $y = x_1 \dots x_j$  and  $z = x_{j+1} \dots x_{2n}$  that are both balanced. By induction,  $y, z \in L(\mathcal{G})$ , and as such  $S \Rightarrow^* y$  and  $S \Rightarrow^* z$ . This implies that

$$S \Rightarrow SS \Rightarrow^* yz = x.$$

Namely,  $x \in L(\mathcal{G})$ .

The remaining case is that  $\alpha_j \neq 0$  for  $j = 2, \dots, 2n - 1$ . If  $x_1 = \mathbf{a}$  then  $\alpha_1 = 1$ . As such, for all  $j = 1, \dots, 2n - 1$ , we must have that  $\alpha_j > 0$ . But then  $\alpha_{2n} = 0$ , which implies that  $\alpha_{2n-1} = 1$ . We conclude that  $x_1 = \mathbf{a}$  and  $x_{2n} = \mathbf{b}$ . As such,  $x_2 \dots x_{2n-1}$  is a balanced word, which by induction is generated by  $L(\mathcal{G})$ . Thus, the  $x$  can be derived via  $S \rightarrow \mathbf{a}S\mathbf{b} \Rightarrow^* \mathbf{a}x_2x_3 \dots x_{2n-1}\mathbf{b} = x$ . Thus,  $x \in L(\mathcal{G})$ .

The case  $x_1 = \mathbf{b}$  is handled in a similar fashion, and implies that  $x \in L(\mathcal{G})$  also in this case. We conclude that  $L_{\mathbf{a}=\mathbf{b}}^{\text{mix}} \subseteq L(\mathcal{G})$ .

Thus  $L_{\mathbf{a}=\mathbf{b}}^{\text{mix}} = L(\mathcal{G})$ . ■

### 1.3 Unbalanced pair

Consider the following language:

$$L_{\mathbf{a} \neq \mathbf{b}} = \left\{ \mathbf{a}^n \mathbf{b}^m \mid n \neq m \text{ and } n, m \geq 0 \right\}.$$

If  $n \neq m$  then either  $n > m$  or  $m > n$ , therefore we can design this grammar by first starting with the basic grammar for when  $n = m$ , and then transition into making more  $\mathbf{a}$ 's or  $\mathbf{b}$ 's.

Let  $X$  be the non terminal representing "choosing" to generate more  $\mathbf{a}$ 's than  $\mathbf{b}$ 's and  $Y$  be the non-terminal for the other case. One grammar that generates  $L_{\mathbf{a} \neq \mathbf{b}}$  will therefore be:

$$S \rightarrow \mathbf{a}S\mathbf{b} \mid \mathbf{a}A \mid \mathbf{b}B, \quad A \rightarrow \mathbf{a}A \mid \epsilon, \quad B \rightarrow \mathbf{b}B \mid \epsilon.$$

### 1.4 Balanced pair in a triple

Consider the language

$$L_4 = \left\{ \mathbf{a}^i \mathbf{b}^j \mathbf{c}^k \mid i = j \text{ or } j = k \right\}.$$

We can essentially combine two copies of the previous grammar (with one version that works on  $\mathbf{b}$  and  $\mathbf{c}$ ) in order to create a grammar that generates  $L_2$ :

$$S \rightarrow S_{\mathbf{a}=\mathbf{b}}C \mid AS_{\mathbf{b}=\mathbf{c}}$$

$$S_{\mathbf{a}=\mathbf{b}} \rightarrow \mathbf{a}S_{\mathbf{a}=\mathbf{b}}\mathbf{b} \mid \epsilon, \quad S_{\mathbf{b}=\mathbf{c}} \rightarrow \mathbf{b}S_{\mathbf{b}=\mathbf{c}}\mathbf{c} \mid \epsilon, \quad A \rightarrow \mathbf{a}A \mid \epsilon, \quad C \rightarrow \mathbf{c}C \mid \epsilon.$$

**Exercise 1.2** Derive a CFG for the language  $L'_4 = \left\{ \mathbf{a}^i \mathbf{b}^j \mathbf{c}^k \mid i = j \text{ or } j = k \text{ or } i = k \right\}$ .

## 1.5 Unbalanced pair in a triple

Now consider the related language

$$L_2 = \left\{ a^i b^j c^k \mid i \neq j \text{ or } j \neq k \right\}.$$

We can essentially combine two copies of the previous grammar (with one version that works on b and c) in order to create a grammar that generates  $L_2$ :

$$\begin{aligned} S &\rightarrow S_{a \neq b} C \mid A S_{b \neq c} & S_{a \neq b} &\rightarrow a S_{a \neq b} b \mid aA \mid bB. & S_{b \neq c} &\rightarrow b S_{b \neq c} c \mid bB \mid cC. \\ A &\rightarrow Aa \mid \epsilon & B &\rightarrow Bb \mid \epsilon & C &\rightarrow Cc \mid \epsilon. \end{aligned}$$

## 1.6 Anything but balanced

Let  $\Sigma = \{a, b\}$ , and let  $\overline{L_{a=b}} = \Sigma^* \setminus \{a^n b^n \mid n \geq 1\}$ .

The idea is that lets first generate all words that contain b in them, and then later the contain a. The grammar for this language is

$$S_1 \rightarrow ZbZaZ \quad Z \rightarrow aZ \mid bZ \mid \epsilon.$$

Clearly  $L(Z) \subseteq \overline{L_{a=b}}$ . The only words we miss, must have all their as before their bs. But these are all words of the form  $a^i b^j$ , where  $i \neq j \geq 0$ . But we already saw how to generate such words in Section 1.3. Putting everything together, we get the following grammar.

$$\begin{aligned} \Rightarrow S &\rightarrow S_1 \mid S_{a \neq b} \\ S_1 &\rightarrow ZbZaZ \\ Z &\rightarrow aZ \mid bZ \mid \epsilon \\ S_{a \neq b} &\rightarrow a S_{a \neq b} b \mid aA \mid bB, \\ A &\rightarrow aA \mid \epsilon, \\ B &\rightarrow bB \mid \epsilon. \end{aligned}$$

## 2 Similar count

Consider the language

$$L = \left\{ w 0^n \mid w \in \{a, b\}^* \text{ and } \#_a(w) = n \right\},$$

where  $\#_a(w)$  is the number of appearances of the character a in  $w$ . The grammar for this language is

$$S \rightarrow \epsilon \mid bS \mid aS0.$$

### 3 Inherent Ambiguity

In lecture, the following ambiguous grammar representing basic mathematical statements was discussed:

$$E \rightarrow E * E \mid E + E \quad N \rightarrow 0N \mid 1N \mid 0 \mid 1.$$

The ambiguity caused because there is no inherent preference from combining expressions with  $*$  over  $+$  or vice versa. It was then fixed by introducing a preference :

$$E \rightarrow E * E \mid T, \quad T \rightarrow N * T \mid N \quad N \rightarrow 0N \mid 1N \mid 0 \mid 1.$$

However some languages are inherently ambiguous, no context free grammar without ambiguity can generate it.

Consider the following language:

$$L = \left\{ a^n b^n c^k d^k \mid n, k \geq 1 \right\} \cup \left\{ a^n b^k c^k d^n \mid n, k \geq 1 \right\}.$$

In other words, it is the language of  $a^+b^+c^+d^+$  where either:

1. the number of a's equals the number of b's and the number c's equals the number of d's
2. the number of a's equals the number of d's and the number of b's equals the number of c's

One ambiguous grammar that generates it

$$S \rightarrow XY \mid Z, \quad X \rightarrow aXb \mid \epsilon, \quad Y \rightarrow cYd \mid \epsilon, \quad Z \rightarrow aZd \mid T, \quad T \rightarrow bTc \mid \epsilon.$$

The reason why all grammars for this language must be ambiguous can be seen in strings of the form  $\left\{ a^n b^n c^n d^n \mid n \geq 1 \right\}$ . Any grammar needs some way of generating the string in a way that either the a's and b's are equal and the c and d's are equal or the a's and d's are equal and the b's and c's are equal. When generating equal a's and b's, it must be still possible to have the same number of c's and d's. When generating equal a's and d's, it must still be possible to have the same number of b's and c's. No matter what grammar is designed, any string of the form  $\left\{ a^n b^n c^n d^n \mid n \geq 1 \right\}$  must have at least two possible parse trees.

(This is of course only an intuitive explanation. A formal proof that any grammar for this language must be ambiguous is considerably more tedious and harder.)

### 4 A harder example

Consider the following language:

$$L = \left\{ xy \mid x, y \in \{0, 1\}^* \text{ where } |x| = |y|, \text{ and } x \neq y \right\}.$$

It should be clear that this language cannot be regular. However, it may not be obvious that we can in fact design a context free grammar for it.  $x$  and  $y$  are guaranteed to be different if, for some  $k$ , the  $k$ th character is 0 in  $x$  and 1 in  $y$  (or vice versa). It is important to notice that we should not try to build  $x$  and  $y$  separately as, in a CFG, we would have no way to enforce them being of the same length. Instead, we just remember that if the string is of length  $2n$ , the first  $n$  characters are considered  $x$  and the second  $n$  characters are  $y$ . Similarly, notice that we cannot choose  $k$  ahead of time for similar reasons.

So, consider the following string

$$w = x_1x_2 \dots x_{k-1} \boxed{1} x_{k+1} \dots x_n y_1y_2 \dots y_{k-1} \boxed{0} y_{k+1} \dots y_n \in L.$$

Then we can rewrite this string as follows

$$w = \overbrace{x_1x_2 \dots x_{k-1}}^{k-1 \text{ chars}} \boxed{1} \overbrace{x_{k+1} \dots x_n}^{n-k \text{ chars}} \overbrace{y_1y_2 \dots y_{k-1}}^{k-1 \text{ chars}} \boxed{0} \overbrace{y_{k+1} \dots y_n}^{n-k \text{ chars}}.$$

In particular, let  $z_1z_2 \dots z_{n-1} = x_{k+1} \dots x_n y_1y_2 \dots y_{k-1}$ . Then,

$$\begin{aligned} w &= \overbrace{x_1x_2 \dots x_{k-1}}^{k-1 \text{ chars}} \boxed{1} z_1z_2 \dots z_{n-1} \boxed{0} \overbrace{y_{k+1} \dots y_n}^{n-k \text{ chars}} \\ &= \underbrace{\overbrace{x_1x_2 \dots x_{k-1}}^{k-1 \text{ chars}} \boxed{1} \overbrace{z_1 \dots z_{k-1}}^{k-1 \text{ chars}}}_{=X} \underbrace{\overbrace{z_k \dots z_{n-1}}^{n-1-k+1 \text{ chars}} \boxed{0} \overbrace{y_{k+1} \dots y_n}^{n-k \text{ chars}}}_{=Y}. \end{aligned}$$

Now,  $X$  is a word of odd length with 1 in the middle (and we definitely know how to generate this kind of words using context free grammars). And  $Y$  is a word of odd length, with 0 in the middle. In particular, any word of  $L$  can be written as either  $XY$  or  $YX$ , where  $X$  and  $Y$  are as above. We conclude, that the grammar for this language is

$$S \rightarrow XY \mid YX \quad X \rightarrow DXD \mid 1 \quad Y \rightarrow DYD \mid 0 \quad D \rightarrow 0 \mid 1.$$