

1. (30 points): (Question 2 from Kozen Homework 5)

Prove that the CFG:

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

generates the set of all strings of  $\{a, b\}$  with equally many  $a$ 's and  $b$ 's. (**Hint:** Characterize elements of the set in terms of the graph of the function  $\#b(y) - \#a(y)$  as  $y$  ranges over prefixes of  $x$ , as we did with balanced parentheses.)

**Solution:**

Let  $A$  be the language of all strings over  $\{a, b\}^*$  that have an equal number of  $a$ 's and  $b$ 's. I'll prove that  $L(G) \subseteq A$  and  $A \subseteq L(G)$  separately.

$L(G) \subseteq A$ :

Let  $w \in L(G)$ . My proof is by induction on the number of steps,  $n$ , in the derivation of  $w$ .

Induction Hypothesis:  $S \xrightarrow{G}^n \alpha \Rightarrow (\#a(\alpha) = \#b(\alpha))$

Base case —  $n = 0$ :

$$\begin{aligned} S &\xrightarrow{G}^0 \alpha \\ \Rightarrow \alpha &= S \\ &(\#a(\alpha) = 0) \wedge (\#b(\alpha) = 0) \\ &\#a(\alpha) = \#b(\alpha) \end{aligned}$$

Induction step —  $S \xrightarrow{G}^n \alpha A \beta \xrightarrow{G}^1 \alpha \mu \beta$  There is a separate case for each production of  $A \rightarrow \mu$  of  $G$ .

$A \rightarrow \mu \equiv S \rightarrow aSb$ :

$$\begin{aligned} \#a(\alpha \mu \beta) &= \#a(\alpha aSb \beta), && \text{case hypothesis: } \mu = aSb \\ &= \#a(\alpha \beta) + \#a(aSb), && \text{properties of addition} \\ &= \#a(\alpha S \beta) + 1, && \#a(s) = 0, \#a(aSb) = 1 \end{aligned}$$

Likewise,  $\#b(\alpha aSb \beta) = \#b(\alpha S \beta) + 1$ .  $\#a(\alpha S \beta) = \#b(\alpha S \beta)$  from which we conclude  $\#a(\alpha aSb \beta) = \#b(\alpha aSb \beta)$  as required.

$A \rightarrow \mu \equiv S \rightarrow aSb$ :

The same argument as for the previous case (swapping  $a$  and  $b$ ) applies here.

$A \rightarrow \mu \equiv S \rightarrow SS$ :

As there are neither any  $a$  nor  $b$  terminals in  $S$ , we have  $\#a(\alpha S \beta) = \#a(\alpha SS \beta)$ , and likewise for the  $b$ 's. Thus,  $\#a(\alpha SS \beta) = \#b(\alpha SS \beta)$  follows directly from the induction hypothesis.

Thus, if  $S \xrightarrow{G}^* w$ , then  $\#a(w) = \#b(w)$ . Therefore,  $w \in A$  and I conclude  $L(G) \subseteq A$ .

$A \subseteq L(G)$ : Let  $w \in A$ . I'll prove  $w \in L(G)$  by induction on  $w$ . As in the proof for the balanced parenthesis language, my proof uses "strong" induction – I use apply the induction hypothesis to strings that are shorter than  $|w| - 1$ .

Induction Hypothesis —  $((w \in A) \wedge (|w| \leq n)) \Rightarrow (w \in L(G))$ :

Base step —  $w = \epsilon$ :

$S \rightarrow \epsilon$ . Thus,  $\epsilon \in L(G)$ .

Induction step:

I break the proof into two cases as was done for the balanced parentheses language according to whether or not  $w$  can be divided into a non-empty prefix and suffix where each has an equal number of  $a$ 's and  $b$ 's.

case  $\exists x, y \in \Sigma^+ . (w = xy) \wedge (\#a(x) = \#b(x)) \wedge (\#a(y) = \#b(y))$ :

Because  $|x| < |w|$  and  $x \in A$ , the induction hypothesis applies, and I conclude  $x \in L(G)$ . Equivalently,  $S \xrightarrow{*}_G x$ . Likewise,  $S \xrightarrow{*}_G y$ . The production  $S \rightarrow SS$  is in  $G$ . Thus,

$$S \xrightarrow{1}_G S S \xrightarrow{*}_G x S \xrightarrow{*}_G x y = w$$

Thus,  $w \in L(G)$ .

case  $\neg \exists x, y \in \Sigma^+ . (w = xy) \wedge (\#a(x) = \#b(x)) \wedge (\#a(y) = \#b(y))$ :

Because  $\#a(x)$  and  $\#b(x)$  are integers and change by 0 or +1 as symbols are appended to  $x$ , I conclude:

$$\forall x, y \in \Sigma^+ . (w = xy) \Rightarrow ((\#a(x) > \#b(x)) \vee (\#a(y) > \#b(y)))$$

I'll assume that  $\#a(x) > \#b(x)$  for all prefixes of  $w$  with lengths in  $\{1 \dots |w| - 1\}$ . The other case is similar. Let  $c_1$  be the first symbol of  $w$  and  $c_{|w|}$  be the last symbol of  $w$ , and choose  $y$  such that  $w = c_1 \cdot y \cdot c_{|w|}$ . By the assumption about  $x$ ,  $\text{count}(c_1) > \#b(c_1)$  which implies that  $c_1 = a$ . Furthermore,  $\#a(c_1 y) > \#b(c_1 y)$ , and  $\#a(c_1 y c_{|w|}) = \#b(c_1 y c_{|w|})$ . Thus,  $c_{|w|} = b$ . In other words,  $w = a \cdot y \cdot b$ . Because

$$\#a(y) = (\#a(w) - 1) = (\#a(w) - 1) = (\#b(w) - 1) = \#b(y)$$

$y \in A$ . By the induction hypothesis,  $y \in L(G)$ . Thus,  $S \xrightarrow{*}_G y$ ; furthermore  $S \xrightarrow{1}_G a S b \xrightarrow{*}_G a y b = w$ . Thus  $w \in L(G)$  as required.

This completes the induction proof.

This induction argument shows that  $A \subseteq L(G)$ .

Having shown  $L(G) \subseteq A$  and  $A \subseteq L(G)$ , I conclude  $L(G) = A$ .

## 2. (30 points): (Question 2 from Kozen Homework 6)

Construct a pushdown automaton that accepts the set of strings in  $\{a, b\}^*$  with equally many  $a$ 's and  $b$ 's. Specify all transitions.

### Solution:

Let  $M$  be a PDA that accepts on empty stack with

$$\begin{aligned} M &= (\{q\}, \{a, b\}, \{\perp, a, b\}, \delta, q, \perp, \emptyset) \\ \delta &= \{ \\ &\quad (q, a, \perp) \rightarrow (q, a \perp) \\ &\quad (q, a, a) \rightarrow (q, aa) \\ &\quad (q, a, b) \rightarrow (q, \epsilon) \\ &\quad (q, b, \perp) \rightarrow (q, b \perp) \\ &\quad (q, b, b) \rightarrow (q, bb) \\ &\quad (q, b, a) \rightarrow (q, \epsilon) \\ &\quad (q, \epsilon, \perp) \rightarrow (q, \epsilon) \\ &\quad \} \end{aligned}$$

This machine has the charming property that

$$\begin{aligned} &(q, xy, \perp) \rightarrow (q, y, \alpha) \\ \Leftrightarrow &((\#a(x) - \#b(x)) = (\#a(\alpha) - \#b(\alpha))) \wedge ((\#a(\alpha) = 0) \vee (\#b(\alpha) = 0)) \end{aligned}$$

This can be shown by induction on  $x$ . Furthermore, the machine can always consume its entire input string, as there is always a move on either input symbol that doesn't empty the stack. Thus, if  $w$  has an equal number of  $a$ 's and  $b$ 's,  $M$  will reach a configuration where  $\#a(\alpha) = \#b(\alpha)$ . Since at least one of  $\#a(\alpha)$  or  $\#b(\alpha)$  must be zero, they must both be zero. Thus, the stack must be either  $\perp$  or  $\epsilon$ . In the former case,

the machine can perform the transition  $(q, \epsilon, \perp) \rightarrow (q, \epsilon)$  to empty its stack. In the other case, the stack is already empty. Thus,  $M$  accepts  $w$ .

Conversely, if  $M$  accepts  $w$ ,  $M$  must read all of  $w$  and empty its stack. By the property shown above, if  $M$  empties its stack, then it has read an equal number of  $a$ 's and  $b$ 's. Thus,  $w$  has an equal number of  $a$ 's and  $b$ 's.

My explanation is longer than is needed to get full credit. My goal is to make sure that everyone in the class gets their questions answered.

3. **(40 points):** Let  $T$  be the language over the alphabet  $\{[, \# , ]\}$  such that every  $[$  is followed by its matching  $\#$ , and every  $\#$  is followed by its matching  $]$ , and the total number of  $[$  symbols and the total number of  $]$  symbols in the string are the same. More formally, let

$$\begin{aligned} \text{left}(\epsilon) &= 0, & \text{middle}(\epsilon) &= 0, & \text{right}(\epsilon) &= 0, \\ \text{left}(x[) &= \text{left}(x) + 1, & \text{middle}(x[) &= \text{middle}(x), & \text{right}(x[) &= \text{right}(x), \\ \text{left}(x\#) &= \text{left}(x), & \text{middle}(x\#) &= \text{middle}(x) + 1, & \text{right}(x\#) &= \text{right}(x), \\ \text{left}(x]) &= \text{left}(x), & \text{middle}(x]) &= \text{middle}(x), & \text{right}(x]) &= \text{right}(x) + 1 \end{aligned}$$

A string  $x$  is in  $T$  iff

$$\forall y, z. x = yz. (\text{left}(y) \geq \text{middle}(y)) \wedge (\text{middle}(y) \geq \text{right}(y)) \wedge (\text{left}(x) = \text{middle}(x) = \text{right}(x))$$

For example, the strings

$$[\#] \quad [[\#[\#]]][\#] \quad [\#][\#]$$

are in  $T$ , but the strings

$$[][\#] \quad [] \quad [\#[\#][\#][\#]]$$

are not.

- (a) **(20 points):** Prove that  $T$  is not a context-free language.

**Solution:**

If  $T$  were context-free, it would have a pumping lemma constant,  $k$ . Let  $z = [^k \#^k ]^k$ . Clearly,  $w \in T$ . Consider any strings  $u, v, w, x, y$  with  $uwxv = z$ ,  $|vx| \geq 1$ , and  $|vwx| \leq k$ . If  $vwx$  is contained in the prefix  $[^k \#^k$ , then pumping it will result in having a different number of  $[$  or  $\#$  symbols than  $]$  symbols. Likewise, if  $vwx$  is a substring of  $\#^k ]^k$ , then pumping it will produce a string with a different number of  $\#$  or  $]$  symbols than  $[$  symbols. In either case, a string that is not in  $T$  is produced. It is not possible for  $vwx$  to contain symbols in both  $[^k$  and  $]^k$  because it would have to have a length of at least  $k + 1$ . Thus,  $T$  doesn't satisfy the conditions of the pumping lemma. Therefore, it is not context-free.

- (b) **(15 points):** Give the grammars for two CFLs,  $A_1$  and  $A_2$  such that  $T = A_1 \cap A_2$ .

**Solution:**

Let

$$\begin{aligned} G_1 &= (\{S_1, R\}, \{[, \# , ]\}, P_1, S_1), & \text{grammar for } A_1 \\ P_1 &= \{ S_1 \rightarrow RS_1R \mid [S_1] \mid S_1S_1 \mid \epsilon \\ & \quad R \rightarrow ] \mid \epsilon \\ & \quad \} \\ G_2 &= (\{S_2, L\}, \{[, \# , ]\}, P_2, S_2), & \text{grammar for } A_2 \\ P_2 &= \{ S_2 \rightarrow LS_2L \mid [S_2] \mid S_2S_2 \mid \epsilon \\ & \quad L \rightarrow [ \mid \epsilon \\ & \quad \} \end{aligned}$$

The grammar  $G_1$  accepts any strings where left and middle parentheses match properly, and allows right parentheses to appear anywhere. The grammar  $G_2$  accepts any strings where middle and right

parentheses match properly, and allows left parentheses to appear anywhere. Thus, any string that is in both  $L(G_1)$  and  $L(G_2)$  has every left parentheses matched by a subsequent middle parentheses, and every middle parentheses matched by a subsequent right parentheses. This is the language  $T$ .

- (c) **(5 points):** Are context-free languages closed under intersection? Give a *short* justification for your answer.

**Solution:**

**No.** As shown above,  $G_1$  and  $G_2$  are context-free grammars.  $L(G_1) \cap L(G_2) = T$ , and  $T$  is not context-free.

- (d) **(5 points):** Are context-free languages closed under complement? Give a *short* justification for your answer.

**Solution:**

As noted in the newsgroup, I had meant to write “intersection” instead of “complement”. A short answer is that CFLs are not closed under complement as shown in the text book. For example, the language  $\{x \mid x = ww\}$  is not a CFL, but its complement is.

Here’s another proof that uses what we’ve just shown. CFLs are closed under union. Let  $A_1$  and  $A_2$  be two CFLs over the same alphabet. Let  $G_1 = (N_1, \Sigma, P_1, S_1)$  and  $G_2 = (N_2, \Sigma, P_2, S_2)$  be CFGs for  $A_1$  and  $A_2$  respectively. We can assume that  $N_1$  and  $N_2$  are disjoint and that neither contains the symbol  $S$  (this can be achieved by renaming non-terminals in one set or the other if needed). Let

$$\begin{aligned} G &= (N, \Sigma, P, S) \\ N &= N_1 \cup N_2 \cup \{S\} \\ P &= P_1 \cup P_2 \cup \{S \rightarrow S_1 \mid S_2\} \end{aligned}$$

It is straightforward to show that  $L(G) = A_1 \cup A_2$ . Thus, CFLs are closed under union.

Now, if CFLs were closed under both union and complement, they would also be closed under intersection by De Morgan’s Law:

$$A_1 \cap A_2 = \sim((\sim A_1) \cup (\sim A_2))$$

We have shown that CFLs are closed under union but not closed under intersection. Therefore, they are not closed under complement.