# Assignment 2, 2019, feedback to Challenges 1–3

For feedback to Challenges 4–6, go to Grok.

## Challenge 1

The main problems were grammars that did not generate the right language.

- Some were missing strings. For example, the grammar:

$$
\begin{aligned}
S &\rightarrow \text{ a } A \text{ a} \mid \text{b } B \text{ b} \\
A &\rightarrow C\ A\ C \\
B &\rightarrow C\ B\ C \\
C &\rightarrow \text{ a} \mid \text{b}
\end{aligned}
$$

  does not generate all the required strings, like a or bbb.

- Others were too generous. For example, the grammar:

$$
\begin{aligned}
S &\rightarrow \text{ a } X \text{ b } X \mid X \text{ b } X \text{ a} \\
X &\rightarrow \text{ a } X \mid \text{a} \mid \epsilon
\end{aligned}
$$

  generates strings like aba: $S \Rightarrow \text{a}Xb X \Rightarrow \text{aba}$.

- Some students still assume that different occurrences of a grammar variable on the RHS of a production must be rewritten to the same string. That, of course, is not the case.

For Challenge 1a, the easiest way to solve it is to distinguish, at the top level of the grammar, the two different cases: Either a is in the middle, or b is in the middle:
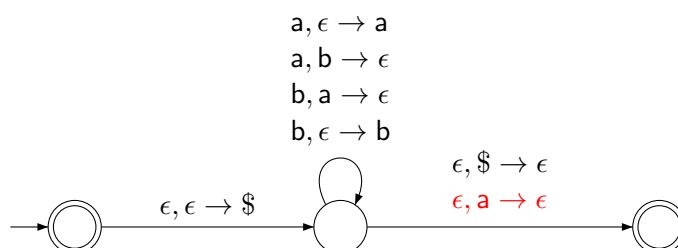
$$
\begin{aligned}
S &\rightarrow \text{ a} \mid \text{b} \mid \text{a } A \text{ a} \mid \text{b } B \text{ b} \\
A &\rightarrow \text{ a} \mid C\ A\ C \\
B &\rightarrow \text{ b} \mid C\ B\ C \\
C &\rightarrow \text{ a} \mid \text{b}
\end{aligned}
$$

Challenge 1b is quite similar to tutorial exercise 80, once we realise that $i \neq j$ is equivalent to $i > j \lor j > i$:

$$
\begin{aligned}
S &\rightarrow T \mid U \\
T &\rightarrow A\ B \\
U &\rightarrow B\ A \\
A &\rightarrow \text{ a} \mid \text{a } A \\
B &\rightarrow \text{ b} \mid \text{a } B \text{ a}
\end{aligned}
$$

## Challenge 2a

Many students did not understand the way a push-down automaton works, and/or they misinterpret '$\epsilon$'. This leads to proposed solutions like this one:

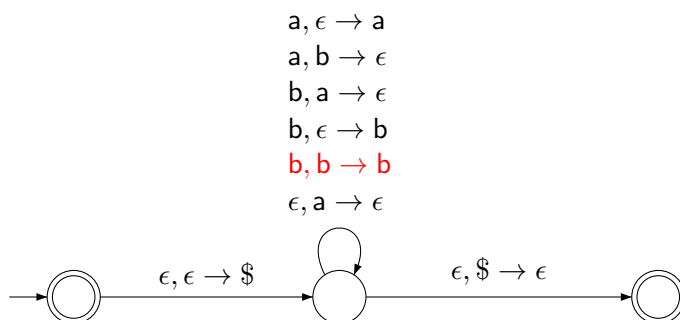This PDF will accept bba, which should be rejected. Here is what the PDA will do, step by step:

a. Consume nothing from the input, move from state 0 to 1, push a $ onto the stack;

b. Consume a b from the input, stay in state 1, push a b to the stack;

c. Consume a b from the input, stay in state 1, push a b to the stack;

d. Consume an a from the input, stay in state 1, push an a to the stack;

e. Consume nothing from the input, transit from state 1 to 2, pop an a from the stack;

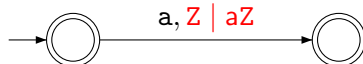f. Since all input has been consumed and the PDA is now in an accept state, it will accept.

Some students reacted to the marking of this type of machine by saying, "But that's not how my PDA is meant to operate. In step d it was obviously meant to pop a b, rather than push an a, and so it would end up rejecting." It is really important to understand how the nondeterminism works. The given PDA will accept bba *because it can*. It is biased towards acceptance.

In a rule like $a, \epsilon \to a$, $\epsilon$ does not stand for "the bottom of the stack". The rule says "consume, if you want, an a from input, which replacing *nothing* on the stack by an a; or, in other words, push a. Nor is $\epsilon$ used as a marker of "end of input".

Another example of a misconception: The red transition allows the following PDA to consume multiple bs from the input for free, as long as the top of the stack is a b:

$$a, \epsilon \to a$$
$$a, b \to \epsilon$$
$$b, a \to \epsilon$$
$$b, \epsilon \to b$$
$$b, b \to b$$
$$\epsilon, a \to \epsilon$$

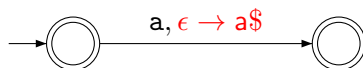$$\epsilon, \epsilon \to \$ \qquad \epsilon, \$ \to \epsilon$$

Some were using unconventional representations without a clear definition or explanation. For example, we cannot possibly know how to read this:
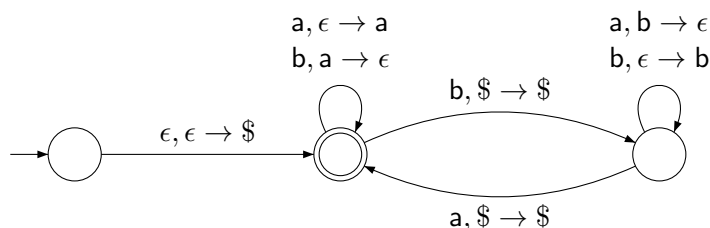
$$a, Z \mid aZ$$

Precision and clarity in the use of formal notation is not optional; you cannot just invent your own notation and hope the reader will guess the intended meaning.

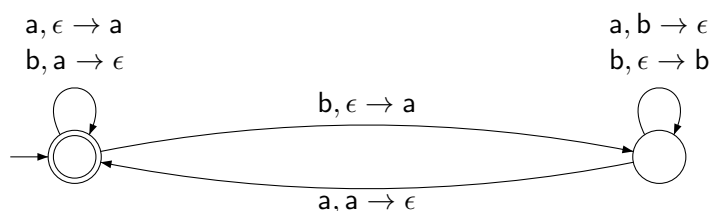Many students resorted to "transitions" that would push or pop several symbols in one go, like this:

$$a, \epsilon \to a\$$$

The way we have defined a PDA, it cannot do that.

Here is a correct 3-state PDA which recognises $C$:

$$a, \epsilon \to a$$
$$b, a \to \epsilon$$
$$b, \$ \to \$$$
$$a, b \to \epsilon$$
$$b, \epsilon \to b$$
$$\epsilon, \epsilon \to \$$$
$$a, \$ \to \$$$

This was the solution we had in mind. I was amazed to see that a 2-state PDA can do the same job. A couple of students found this solution, avoiding the use of $:

$$a, \epsilon \to a$$
$$b, a \to \epsilon$$
$$b, \epsilon \to a$$
$$a, b \to \epsilon$$
$$b, \epsilon \to b$$
$$a, a \to \epsilon$$

Brilliant.

## Challenge 2b

The most common comments from the markers were:

- Missing a as a base case.

- Informal proof. For example, clear statement of the assumption (that is, the immediate substructure is "a-leaning") in each induction step.

- Inexplicit induction hypotheses.

- Not proving in both directions.

Here is an example of a valid proof:

Say that a string $w \in \{a, b\}^*$ is *a-leaning* if it contains at least as many as as bs. We show that $G$ generates the language of a-leaning strings.

(1) That every string in $L(G)$ is a-leaning can be shown by structural induction. Clearly $\epsilon$ and a are a-leaning. For the induction step, note that if $w$ and $w'$ are a-leaning, so are $awb$, $bwa$, and $ww'$.

(2) That every a-leaning string is in $L(G)$ can be shown by induction on the length of strings. Consider any a-leaning $w$. If $|w| \leq 1$ then $w = \epsilon$ or $w = a$ and hence $w$ can be generated from $S$. For the induction step, assume that every a-leaning $w$ of length less than $n$ can be generated from $S$. We show that any a-leaning $w$ with $|w| = n$ can be generated from $S$.

Case (2.1): If $w$ starts and ends with different symbols, then $w = aw'b$ or $w = bw'a$ for some a-leaning $w'$. By the induction hypothesis, $w'$ can be generated from $S$. But then so can $w$, owing to either rule $S \to aSb$ or rule $S \to bSa$.

Case (2.2): If $w$ starts and ends with b then some non-empty proper prefix $w_1$ of $w$ must have the same number of as and bs. Hence $w = w_1w_2$ with $w_1 \neq \epsilon$, $w_2 \neq \epsilon$, both of $w_1$ and $w_2$ a-leaning and of length less than $n$. By the induction hypothesis, both can be generated from $S$. But then so can $w$, owing to rule $S \to SS$.

Case (2.3): Otherwise $w$ starts and ends with a. If $w$ has the same number of as and bs then some non-empty proper prefix of $w$ must have the same number of as and bs as well, and the argument then goes as in case (2.2). Otherwise $w$ has strictly more as than bs, in which case $w = aw'$ for some a-leaning $w'$. But a can be generated from $S$, and so can $w'$, by the induction hypothesis. Hence $w$ can be generated from $S$.

In all cases, $w \in L(G)$.

## Challenge 3a

No dramas there. $R^3$ is simply the concatenation of three regular languages. The class of regular languages is closed under concatenation, so $R^3$ is regular.

## Challenge 3b

Here the main issue was sloppy use of the pumping lemma.

$thrice(R)$ is not necessarily regular, and this should not come as a surprise, since we already know that $\{ww \mid w \in \Sigma^*\}$ is not regular (or even context-free). We can prove the claim formally with the pumping lemma for regular languages. For example, consider $R = \{0, 1\}^*$, that is, the set of all finite binary strings. We use the pumping lemma for regular languages to show that $thrice(R)$ is not regular in this case.

Assume $thrice(R)$ is regular. Let $p$ be the pumping length, and consider the string

$$s = 0^p 1 0^p 1 0^p 1$$

Clearly, $s \in thrice(R)$ and $s$ has length greater than $p$. The pumping lemma tells us that $s$ can be split into three segments, $s = xyz$, with $y \neq \epsilon$ and $|xy| \leq p$, in such a way that $xy^i z \in thrice(R)$ for all $i \in \mathbb{N}$. From $|xy| \leq p$ it follows that (the non-empty) $y$ consists of zeros only. Hence $xz$ cannot be of the form required to be in $thrice(R)$. It follows that $thrice(R)$ is not regular.

## Challenge 3c

This was the trickiest challenge up to this point, and quite a few simply skipped this sub-question. There are lots of possible solutions—all we have to do is find a single regular language $R$ and show that $snip(R)$ is not regular. There is no algorithm that I know of that will help find that $R$; it is a matter of trying out some simple (and maybe less simple) candidate $R$s and understand what $snip(R)$ becomes. This challenge is a good example of how we can make use of closure properties to good effect. For example, here is a solution:

Take the regular language $R = a^*bc^*$ over alphabet $\{a, b, c\}$. The language $snip(R)$ contains strings of two different forms:

a. $a^i c^i$, with $i \in \mathbb{N}$.

b. $a^i b c^j$, with $i, j \in \mathbb{N}$ (and $i \neq j$).

The first form comes from those strings in $R$ for which the length is a multiple of 3, and b falls in the middle third of the string. The second form comes from those strings in $R$ for which the length is a multiple of 3, and b falls in the first or last third.

Hence $snip(R) \cap a^* c^* = \{a^i c^i \mid i \in \mathbb{N}\}$. That is, when we intersect $snip(R)$ with the regular language $a^* c^*$, we get a non-regular language. Since the class of regular languages is closed under intersection, this means $snip(R)$ is not regular.