

## Problem Set 7

---

What can you do with regular expressions? What are the limits of regular languages? In this problem set, you'll explore the answers to these questions along with their practical consequences.

As always, please feel free to drop by office hours, ask on Piazza, or send us emails if you have any questions. We'd be happy to help out.

Good luck, and have fun!

**Due Friday, February 28<sup>th</sup> at 2:30PM**

## Problem One: Designing Regular Expressions

Below are a list of alphabets and languages over those alphabets. For each language, write a regular expression for that language. ***Please use our online tool to design, test, and submit your regular expressions. Typed or handwritten solutions will not be accepted.*** To use it, visit the CS103 website and click the “Regex Editor” link under the “Resources” header. If you submit in a pair, please tell us in your GradeScope submission who submitted your answers to this question. Also, as a reminder, please test your submissions thoroughly, since we’ll be grading them with an autograder.

- i. Let  $\Sigma = \{a, b, c, d, e\}$ . Write a regular expression for the language  $\{w \in \Sigma^* \mid \text{the letters in } w \text{ are sorted alphabetically}\}$ .
- ii. Write a regular expression for the complement of the language from part (i) of this problem.

*As with NFAs, there’s no simple way to start with a regex for a language  $L$  and to turn it into a regex for  $\bar{L}$ .*

- iii. On Unix-style operating systems like macOS or Linux, files are organized into directories. You can reference a file by giving a *path* to the file, a series of directory names separated by slashes. For example, the path `/home/username/` might represent a user’s home directory, and a path like `/home/username/Documents/PS7.tex` might represent that person’s solution to this problem set. Paths that start with a slash character are called *absolute paths* and say exactly where the file is on disk. Paths that don’t start with a slash are called *relative paths* and say where, relative to the current folder, a file can be found. For example, if I’m logged into my computer and am in my home folder, I could look up the file `Documents/PS7.tex` to find my solution to this problem set.

The general pattern here is that a file path consists of a series of directory or file names separated by slashes. That path might optionally start with a slash, but isn’t required to, and it might optionally end with a slash, but isn’t required to. However, you can’t have two consecutive slashes.\*

Let  $\Sigma = \{a, /\}$ . Write a regular expression for  $L = \{w \in \Sigma^* \mid w \text{ represents the name of a file path on a Unix-style system}\}$ . For example, `/aaa/a/aa`  $\in L$ , `/`  $\in L$ , `a`  $\in L$ , `/a/a/a/`  $\in L$ , and `aaa/`  $\in L$ , but `//a//`  $\notin L$ , `a//a`  $\notin L$ , and  $\epsilon \notin L$ .

Fun fact: this problem comes from former TA Amy Liu, who fixed a bug in industrial code that arose when someone wrote the wrong regex for this language. Oops.

- iv. Suppose you are taking a walk with your dog on a leash of length two. Let  $\Sigma = \{y, d\}$  and let  $L = \{w \in \Sigma^* \mid w \text{ represents a walk with your dog on a leash where you and your dog both end up at the same location}\}$ . For example, we have `yyddddyy`  $\in L$  because you and your dog are never more than two steps apart and both of you end up four steps ahead of where you started; similarly, `ddyddy`  $\in L$ . However, `yyyddd`  $\notin L$ , since halfway through your walk you’re three steps ahead of your dog; `ddy`  $\notin L$ , because your dog ends up two steps ahead of you; and `ddyddy`  $\notin L$ , because at one point your dog is three steps ahead of you. Write a regular expression for  $L$ .

*Note that, unlike Problem Set Six, you and your dog **must** end at the same position.*

\* In some cases you technically *can* have multiple consecutive slashes, but we’ll ignore that for now.

## Problem Two: Finite Languages

A language  $L$  is called *finite* if  $L$  contains finitely many strings (that is,  $|L|$  is a natural number). It turns out that all finite languages are regular. Given a finite language  $L$ , explain how to write a regular expression for  $L$ . Briefly justify your answer; no formal proof is necessary.

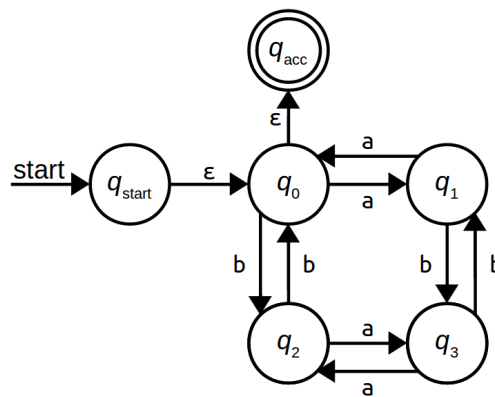
*Watch for edge cases!*

*As an optional thought exercise, you might also think about how you could make an NFA or a DFA for any finite language. As a hint, look up the trie data structure on Wikipedia.*

## Problem Three: State Elimination

The state elimination algorithm gives a way to transform a finite automaton (a DFA or NFA) into a regular expression. It's a really beautiful algorithm once you get the hang of it, so we thought that we'd let you try it out on a particular example.

Let  $\Sigma = \{a, b\}$  and let  $L = \{w \in \Sigma^* \mid w \text{ has an even number of } a\text{'s and an even number of } b\text{'s}\}$ . Below is a finite automaton for  $L$  that we've prepared for the state elimination algorithm by adding in a new start state  $q_{start}$  and a new accept state  $q_{acc}$ :



We'd like you to use the state elimination algorithm to produce a regular expression for  $L$ .

- i. Run two steps of the state elimination algorithm on the above automaton. Specifically, first remove state  $q_1$ , then remove state  $q_2$ . Show your result at this point.

*Go slowly. Remember that to eliminate a state  $q$ , you should identify all pairs of states  $q_{in}$  and  $q_{out}$  where there's a transition from  $q_{in}$  to  $q$  and from  $q$  to  $q_{out}$ , then add shortcut edges from  $q_{in}$  to  $q_{out}$  to bypass state  $q$ . Remember that  $q_{in}$  and  $q_{out}$  can be the same state. If you've done everything right at the end of this stage, none of the transitions you have at this point should have Kleene stars in them.*

- ii. Finish the state elimination algorithm, showing your work. **Submit your resulting regular expression through our regular expression tool** along the lines of what you did in Prob. 1.

*You should start seeing Kleene stars appearing as you remove the remaining states. Not sure whether you have the right answer? **Test your result thoroughly!***

- iii. Without making reference to the original automaton given above, give an intuitive explanation for how the regular expression you found in part (ii) works.

*The regex you've found works by matching strings in a very creative way. Tell us what that way is.*



although it's a string of balanced curly braces, the nesting goes five levels deep.

- ii. Look back at your proof from part (i) of this problem. Imagine that you were to take that exact proof and blindly replace every instance of “ $L_1$ ” with “ $L_2$ .” This would give you a (incorrect) proof that  $L_2$  is nonregular (which we know has to be wrong because  $L_2$  is indeed regular—if you need convincing think about how you might design a DFA for it). Where would the error be in that proof? Be as specific as possible.

*Again, you should be able to point at a specific spot in the proof that contains a logic error and explain exactly why the statement in question is not true or not supported by the preceding statements. If you can't do this, it likely means you have an error in your proof from part (i)!*

## Problem Six: State Lower Bounds

The Myhill-Nerode theorem we proved in lecture is actually a special case of a more general theorem about regular languages that can be used to prove lower bounds on the number of states necessary to construct a DFA for a given language.

- i. Let  $L$  be a language over  $\Sigma$ . Suppose there's a set  $S$ , which may be finite and which may be infinite, such that any two distinct strings  $x, y \in S$  are distinguishable relative to  $L$  (that is,  $x \not\equiv_L y$  for any two strings  $x, y \in S$  where  $x \neq y$ .) Prove that any DFA for  $L$  must have at least  $|S|$  states. (You sometimes hear this referred to as *lower-bounding* the size of any DFA for  $L$ .)

According to old-school Twitter rules, all tweets need to be 140 characters or less. Let  $\Sigma$  be the alphabet of characters that can legally appear in a tweet and consider the following language:

$$TWEETS = \{ w \in \Sigma^* \mid |w| \leq 140 \}.$$

This is the language of all legal tweets, assuming the empty string is a legal tweet. The good news is that this language is regular. The bad news is that any DFA for it has to be pretty large.

- ii. Using your result from part (i), prove that any DFA for  $TWEETS$  must have at least 142 states.

*It might be easier to tackle this problem if you consider replacing 140 and 142 with some smaller numbers (say, 2 and 4) to build up an intuition. And work backwards – what will you need to do to invoke part (i)?*

- iii. Define a 142-state DFA for  $TWEETS$  using the formal 5-tuple definition of a DFA. Briefly explain how your DFA works. No formal proof is necessary.

*Again, this might be a lot easier to do if you first reduce 140 and 142 to 2 and 4, respectively, and see what you come up with. Start by drawing out what the DFA would look like, then think about how you'd formalize your idea as a 5-tuple. Remember that to define the transition function, you'll want to use our methods of formally defining a function that we learned earlier in the course; in particular look at how you defined piece-wise functions for Pset4 Q1.*

Your results from parts (ii) and (iii) show that the smallest possible DFA for  $TWEETS$  has exactly 142 states. This approach to finding the smallest object of some type – using some theorem to prove a lower bound (“we need at least this many states”) combined with a specific object of the given type (“we can't do worse than this”) is a common strategy in algorithm design and computational complexity theory. If you take classes like CS161, CS254, etc., you'll likely see similar sorts of approaches!

## Problem Seven: Regular Languages and Equivalence Relations

Throughout this problem set you've been working with the idea that we can take a language  $L$  over some alphabet  $\Sigma$ , then work with its distinguishability relation  $\equiv_L$ . A closely related binary relation is the *indistinguishability* relation for  $L$ , denoted  $\equiv_L$ . It's also a binary relation over  $\Sigma^*$ , and its definition is the negation of the one for distinguishability:

$$x \equiv_L y \text{ if } \forall w \in \Sigma^*. (xw \in L \leftrightarrow yw \in L).$$

Amazingly, this is always an equivalence relation, regardless of what  $L$  is!

- i. Prove that if  $L$  is a language over  $\Sigma$ , then  $\equiv_L$  is an equivalence relation over  $\Sigma^*$ .

Whenever you see an equivalence relation, you should immediately start thinking about what its equivalence classes are. Doing so will usually tell you something interesting.

Let's make this more concrete. Let  $\Sigma = \{a, b\}$  and consider the language  $M = \{ w \in \Sigma^* \mid \text{the number of } b\text{'s in } w \text{ is congruent to 1 modulo 5 or to 3 modulo 5} \}$ . For example,  $\mathbf{aba} \in M$ ,  $\mathbf{baaabbaaab} \in M$ , and  $\mathbf{bbbbbb} \in M$ , but  $\mathbf{aa} \notin M$  and  $\mathbf{abba} \notin M$ .

- ii. How many equivalence classes does the  $\equiv_M$  relation have? Briefly describe what those equivalence classes are and give a system of representatives for  $\equiv_M$ .

*Need a refresher on systems of representatives? Check out Problem Set Three.*

You might have noticed that each equivalence class of  $\equiv_M$  either consists of a bunch of strings not in  $M$  or of a bunch of strings that are in  $M$ . That's not a coincidence!

- iii. Let  $L$  be a language over some alphabet  $\Sigma$  and let  $x \in \Sigma^*$  be some string. Prove that either *every string* in  $[x]_{\equiv_L}$  is in  $L$  or that *no strings* in  $[x]_{\equiv_L}$  are.

The number of equivalence classes of an equivalence relation is called its *index*; the index of an equivalence relation  $R$  is denoted  $I(R)$ . This quantity might be finite, or it might be an infinite cardinality like  $\aleph_0$ , or even one of the infinities bigger than that.

Armed with the idea of an index, we can state a powerful theorem about finite automata:

**Theorem:** If  $L$  is a language over  $\Sigma$ , then every DFA for  $L$  must have at least  $I(\equiv_L)$  states.

In other words, there's a connection between the number of equivalence classes of a particular binary relation and the minimum sizes of DFAs for that language!

- iv. Prove the above theorem. Feel free to use the *axiom of choice*, which says that every equivalence relation has at least one system of representatives.

*Proving this theorem is mostly an exercise in connecting together ideas you've seen used in other places. Think about the relationship between indices and systems of representatives, between distinguishability and indistinguishability, and between what you're doing here and what you've done earlier on this problem set.*

There's a very nice intuition for what this theorem says. You can think of the indistinguishability relation for a language  $L$  as pinning down the idea "a DFA for  $L$  can't tell the difference between these two strings." If you think back to our intuition behind DFA design – build a DFA where each state keeps track of some different piece of information – then you can think of  $I(\equiv_L)$  as capturing the number of different pieces of information you'd need to remember. The theorem then says that if you want to build a DFA for a language  $L$ , you'll need at least one state per piece of information.