

## 1 Regular Expressions

Give regular expressions for each of the following languages over the alphabet  $\{0, 1\}$ .

1. All strings containing the substring  $000$ .

**Solution:** (*Intuition*) Having a substring means a string  $w$  must contain the substring  $x$ , but also that any other substrings must appear before or after  $x$ .

(*Strategy*) Each string in the set has three parts: any-string + substring + any-string. any-string can be described by  $\Sigma^* = (0 + 1)^*$ . We know the described substring....

(*Result*)  $(0 + 1)^*000(0 + 1)^*$  ■

2. All strings *not* containing the substring  $000$ .

**Solution:** (*Intuition*) For a string to not allow the substring  $w$ , we can manually define our string to only contain patterns that do not contain  $w$ .

(*Strategy*) Since the substring is  $000$  and our alphabet contains of  $0$ s and  $1$ s, the only patterns that are allowed repetitively are  $1$ ,  $01$  or  $001$ :  $(1 + 01 + 001)^*$ . Since the allowed patterns ends with  $1$ , in order to also accommodate strings that can end with  $0$ s (with a maximum of 2  $0$ s at that to avoid our substring  $000$ ) the last substring would be:  $(\epsilon + 0 + 00)$

(*Result*)  $(1 + 01 + 001)^*(\epsilon + 0 + 00)$  ■

**Solution:** (*Intuition*) To not contain the substring  $w$ , with  $n$  consecutive repetitions of the symbol  $a$ , we must ensure that a string  $w$  must only allow substrings containing 0 to  $n - 1$   $a$ 's at any point which can precede or succeed substrings consisting of other symbols in the alphabet.

(*Strategy*) The substring that describes a max of two  $0$ s:  $(\epsilon + 0 + 00)$ . Since the only other symbol in our alphabet is  $1$ , sandwiching it between our previously defined substring, while allowing repetitions, gives us :

(*Result*)  $(\epsilon + 0 + 00)(1(\epsilon + 0 + 00))^*$  ■

3. All strings in which every run of  $0$ s has length at least 3.

**Solution:** (*Intuition*) Our alphabet only has  $0$ s and  $1$ s. Every run of  $0$ s with length of atleast 3 can be preceded or succeeded by any number of  $1$ s.

(*Strategy*) Each string in the set has 2 parts :  $1$ s of any length +  $0$ s of length  $\geq 3$ . Runs of  $0$ s with length of at least 3 can be described as  $0000^*$ .

(*Result*)  $(1 + 0000^*)^*$  ■

**Solution:** (*Intuition*) Since every run of 0s should have length of at least 3, We can breakdown the string to three substrings  $x, y, z$ . We can define the main pattern for the body of the string to allow either 1s or 0s of minimum length 3 in substring  $y$ . Substrings  $x$  and  $z$  are defined to accommodate the beginning and ending patterns respectively.

(*Strategy*) We define  $y$ , the main body of the string, to be alternating occurrences of either any number of 1s or 0s of length greater than 3, which gives us  $((\epsilon + 0000^*)1)^*$ . To accommodate strings that start with 1,  $x$  can be described as  $(\epsilon + 1)$  and for strings that end with three or more zeroes,  $z$  can be described as  $(\epsilon + 0000^*)$

(*Result*)  $(\epsilon + 1)((\epsilon + 0000^*)1)^*(\epsilon + 0000^*)$  ■

4. All strings in which 1 does not appear after a substring 000.

**Solution:** (*Intuition*) We can have a maximum of two 0s before any 1. So we can split the problem to two substrings  $x$  and  $y$  and write it as  $xy$ .  $x$  represents all substrings with 1s where there are less than three 0s before a 1.  $y$  represents substrings containing only 0s.

(*Strategy*) For  $x$ , we can have either no 0 before a 1 or one 0 before a 1 or two 0s before a 1. This can be represented by  $(1 + 01 + 001)^*$ . For  $y$ , we can use  $0^*$  since it represents zero or more 0s.

(*Result*)  $(1 + 01 + 001)^*0^*$  ■

5. All strings containing at least three 0s.

**Solution:** (*Intuition*) Containing at least three 0s can be interpreted as having three substrings say  $x, y, z$ . Then this problem becomes similar to problem 1 in that we can have any other substring before, after and inbetween  $x, y$  and  $z$ . (*Strategy*) Each string has seven parts: any-string + substring 1 + any-string + substring 2 + any-string + substring 3 + any-string. Any-string can be described by  $\Sigma^* = (0 + 1)^*$  and our substrings are just 0.  $(0 + 1)^*0(0 + 1)^*0(0 + 1)^*0(0 + 1)^*$  ■

**Solution (clever):** (*Intuition*) This problem can instead be viewed from the lens of counting the first three 0s. Because our alphabet only have 0s and 1s, this means that the only symbols around the first three 0s can be 1s. Then any-string can follow after the first three 0s. The resulting regular expression is simpler than the previous solution. The second clever solution is the same but in reverse. Reversing the expression works because the substring 0 is symmetric but it can also be viewed as counting the last three 0s. (*Strategy*) Each string has seven parts: 1s + 0 + 1s + 0 + 1s + 0 + any-string. Any-string can be described by  $\Sigma^* = (0 + 1)^*$ , 1s is just  $1^*$  and 0 is just 0. The second clever solution is just this in reverse.  $1^*01^*01^*0(0 + 1)^*$  or  $(0 + 1)^*01^*01^*01^*$  ■

6. Every string except **000**. [Hint: Don't try to be clever.]

**Solution:** Every string  $w \neq 000$  satisfies one of three conditions: Either  $|w| < 3$ , or  $|w| = 3$  and  $w \neq 000$ , or  $|w| > 3$ . The first two cases include only a finite number of strings, so we just list them explicitly. The last case includes *all* strings of length at least 4.

$$\begin{aligned} &\epsilon + 0 + 1 + 00 + 01 + 10 + 11 \\ &+ 001 + 010 + 011 + 100 + 101 + 110 + 111 \\ &+ (1 + 0)(1 + 0)(1 + 0)(1 + 0)(1 + 0)^* \end{aligned}$$

■

**Solution (clever):** A string  $w$  in the language falls in one of the following three categories:

- (a) Does not contain **1** and  $|w| < 3$ .
- (b) Contains **1** in the first three symbols.
- (c) Does not contain **1** in the first three symbols and  $|w| > 3$ .

The strings in (a) can be listed explicitly. The strings in (b) can be further categorized based on the location of the first occurrence of **1**, and can be written as  $(1 + 01 + 001)(1 + 0)^*$ . Notice that strings with multiple **1**s in the first three symbols are also captured by the given expression. For example, **111** can be generated by  $1(1 + 0)^*$ . Finally, the strings in (c) can be written as  $(000(1 + 0))(1 + 0)^*$ . Combining the three categories, we obtain the following expression:

$$\epsilon + 0 + 00 + (1 + 01 + 001 + 000(1 + 0))(1 + 0)^*$$

■

**Solution:** CA Ajitesh Dasaratha wrote an excellent in-depth explainer for this problem. It is on the website: [https://ecealgo.com/materials/extra\\_content/RegEx\\_Ajitesh\\_1.pdf](https://ecealgo.com/materials/extra_content/RegEx_Ajitesh_1.pdf). It does go into topics that we'll explore in a few weeks such as DFAs and DFA to RegEx conversion so keep this material handy. It'll age well!

■

7. All strings  $w$  such that in every prefix of  $w$ , the number of 0s and 1s differ by at most 1.

**Solution:** Equivalently, strings that alternate between 0s and 1s:  $(01 + 10)^*(\epsilon + 0 + 1)$  ■

8. Any string not in  $0^* + 1^*$

**Solution:** So the language represented by  $0^* + 1^*$  contains all the binary strings that have only 0's and all the binary strings that have only 1's. Hence, all the strings not included in that language have at least one 0 and one 1. Something like this would do:

$$(0 + 1)^* 0 (0 + 1)^* 1 (0 + 1)^* + (0 + 1)^* 1 (0 + 1)^* 0 (0 + 1)^*$$

Note that we need to account for 01 and 10 so that's why we have to have two parts to the above expression. ■

- \*9. All strings containing at least two 0s and at least one 1.

**Solution:** There are three possibilities for how such a string can begin:

- Start with 00, then any number of 0s, then 1, then anything.
- Start with 01, then any number of 1s, then 0, then anything.
- Start with 1, then a substring with exactly two 0s, then anything.

All together:  $000^*1(0 + 1)^* + 011^*0(0 + 1)^* + 11^*01^*0(0 + 1)^*$

Or equivalently:  $(000^*1 + 011^*0 + 11^*01^*0)(0 + 1)^*$  ■

**Solution:** There are three possibilities for how the three required symbols are ordered:

- Contains a 1 before two 0s:  $(0 + 1)^* 1 (0 + 1)^* 0 (0 + 1)^* 0 (0 + 1)^*$
- Contains a 1 between two 0s:  $(0 + 1)^* 0 (0 + 1)^* 1 (0 + 1)^* 0 (0 + 1)^*$
- Contains a 1 after two 0s:  $(0 + 1)^* 0 (0 + 1)^* 0 (0 + 1)^* 1 (0 + 1)^*$

So putting these cases together, we get the following:

$$\begin{aligned} & (0 + 1)^* 1 (0 + 1)^* 0 (0 + 1)^* 0 (0 + 1)^* \\ & + (0 + 1)^* 0 (0 + 1)^* 1 (0 + 1)^* 0 (0 + 1)^* \\ & + (0 + 1)^* 0 (0 + 1)^* 0 (0 + 1)^* 1 (0 + 1)^* \end{aligned}$$

■

**Solution (clever):**  $(0 + 1)^* (101^*0 + 011^*0 + 01^*01)(0 + 1)^*$  ■

- ★ 10. All strings in which the substring **000** appears an even number of times.  
(For example, **0001000** and **0000** are in this language, but **00000** is not.)

**Solution:** Every string in  $\{0, 1\}^*$  alternates between (possibly empty) blocks of **0**s and individual **1**s; that is,  $\{0, 1\}^* = (0^*1)^*0^*$ . Trivially, every **000** substring is contained in some block of **0**s. Our strategy is to consider which blocks of **0**s contain an even or odd number of **000** substrings.

Let  $X$  denote the set of all strings in  $0^*$  with an even number of **000** substrings. We easily observe that  $X = \{0^n \mid n = 1 \text{ or } n \text{ is even}\} = 0 + (00)^*$ .

Let  $Y$  denote the set of all strings in  $0^*$  with an odd number of **000** substrings. We easily observe that  $Y = \{0^n \mid n > 1 \text{ and } n \text{ is odd}\} = 000(00)^*$ .

We immediately have  $0^* = X + Y$  and therefore  $\{0, 1\}^* = ((X + Y)1)^*(X + Y)$ .

Finally, let  $L$  denote the set of all strings in  $\{0, 1\}^*$  with an even number of **000** substrings. A string  $w \in \{0, 1\}^*$  is in  $L$  if and only if an even number of blocks of **0**s in  $w$  are in  $Y$ ; the remaining blocks of **0**s are all in  $X$ . We consider two simple subcases to help us make sure we don't make mistakes.

- Exactly zero blocks of **0**s are in  $Y$ . This means all blocks of **0**s are in  $X$  and they alternate with blocks of **1**s. We capture this by  $(X1)^*X$  where the last  $X$  is to allow ending in blocks of **0**s.
- Exactly two blocks of **0**s are in  $Y$ . The expression  $(X1)^*Y1 \cdot (X1)^*Y(\epsilon + 1(X1)^*X)$  captures this where we are paying attention in the expression to what comes after the second block of  $Y$ ; either we end in that block or we need a **1** to separate a string that has zero blocks from  $Y$ .

Now we consider the general case with an even number of blocks in  $Y$  which can be zero, two or a multiple of two that can be simulated by repetitions. We can obtain this by considering the expression for exactly two blocks and allowing concatenation arbitrary number of times while being careful to ensure that we separate with a **1** as needed and allowing for the expression to end appropriately. The expression  $((X1)^*Y1 \cdot (X1)^*Y1)^*$  allows an arbitrary even number of blocks in  $Y$  with the caveat that it always ends with  $Y1$  if it has non-zero repetitions. We need to allow ending in  $Y$  or  $(X1)^*X$  as in the second subcase above. Letting  $Z = (X1)^*Y1 \cdot (X1)^*Y$  helps us simplify the final expression.

Putting the preceding observations together we obtain the following expression.

$$L = (Z1)^*(Z + (X1)^*X)$$

We will not plug in the expressions for  $X$  and  $Y$  because it will make the expression too long. There are likely to be shorter expressions based on better case analysis but debugging them is probably not worth it unless it is very important in some application.

Whew! ■