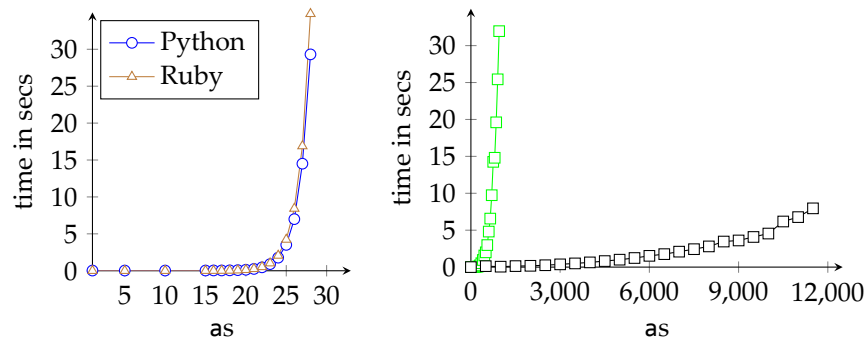


## Handout 2 (Regular Expression Matching)

This lecture is about implementing a more efficient regular expression matcher (the plots on the right)—more efficient than the matchers from regular expression libraries in Ruby and Python (the plots on the left). These plots show the running time for the evil regular expression  $a^{?{n}}a^{\{n\}}$  and string composed of  $n$  as. We will use this regular expression and strings as running example. To see the substantial differences in the two plots below, note the different scales of the  $x$ -axes.



Having specified in the previous lecture what problem our regular expression matcher is supposed to solve, namely for any given regular expression  $r$  and string  $s$  answer *true* if and only if

$$s \in L(r)$$

we can look at an algorithm to solve this problem. Clearly we cannot use the function  $L$  directly for this, because in general the set of strings  $L$  returns is infinite (recall what  $L(a^*)$  is). In such cases there is no way we can implement an exhaustive test for whether a string is member of this set or not. In contrast our matching algorithm will operate on the regular expression  $r$  and string  $s$ , only, which are both finite objects. Before we come to the matching algorithm, however, let us have a closer look at what it means when two regular expressions are equivalent.

### Regular Expression Equivalences

We already defined in Handout 1 what it means for two regular expressions to be equivalent, namely if their meaning is the same language:

$$r_1 \equiv r_2 \stackrel{\text{def}}{=} L(r_1) = L(r_2)$$

It is relatively easy to verify that some concrete equivalences hold, for example

$$\begin{aligned}
(a + b) + c &\equiv a + (b + c) \\
a + a &\equiv a \\
a + b &\equiv b + a \\
(a \cdot b) \cdot c &\equiv a \cdot (b \cdot c) \\
c \cdot (a + b) &\equiv (c \cdot a) + (c \cdot b)
\end{aligned}$$

but also easy to verify that the following regular expressions are *not* equivalent

$$\begin{aligned}
a \cdot a &\not\equiv a \\
a + (b \cdot c) &\not\equiv (a + b) \cdot (a + c)
\end{aligned}$$

I leave it to you to verify these equivalences and non-equivalences. It is also interesting to look at some corner cases involving  $\epsilon$  and  $\emptyset$ :

$$\begin{aligned}
a \cdot \emptyset &\not\equiv a \\
a + \epsilon &\not\equiv a \\
\epsilon &\equiv \emptyset^* \\
\epsilon^* &\equiv \epsilon \\
\emptyset^* &\not\equiv \emptyset
\end{aligned}$$

Again I leave it to you to make sure you agree with these equivalences and non-equivalences.

For our matching algorithm however the following seven equivalences will play an important role:

$$\begin{aligned}
r + \emptyset &\equiv r \\
\emptyset + r &\equiv r \\
r \cdot \epsilon &\equiv r \\
\epsilon \cdot r &\equiv r \\
r \cdot \emptyset &\equiv \emptyset \\
\emptyset \cdot r &\equiv \emptyset \\
r + r &\equiv r
\end{aligned}$$

which always hold no matter what the regular expression  $r$  looks like. The first two are easy to verify since  $L(\emptyset)$  is the empty set. The next two are also easy to verify since  $L(\epsilon) = \{\epsilon\}$  and appending the empty string to every string of another set, leaves the set unchanged. Be careful to fully comprehend the fifth and sixth equivalence: if you concatenate two sets of strings and one is the empty set, then the concatenation will also be the empty set. To see this, check the definition of  $\_ @ \_$ . The last equivalence is again trivial.

What will be important later on is that we can orient these equivalences and read them from left to right. In this way we can view them as *simplification rules*. Consider for example the regular expression

$$(r_1 + \emptyset) \cdot \epsilon + ((\epsilon + r_2) + r_3) \cdot (r_4 \cdot \emptyset) \tag{1}$$

If we can find an equivalent regular expression that is simpler (smaller for example), then this might potentially make our matching algorithm run faster.

The reason is that whether a string  $s$  is in  $L(r)$  or in  $L(r')$  with  $r \equiv r'$  will always give the same answer. In the example above you will see that the regular expression is equivalent to  $r_1$ . You can verify this by iteratively applying the simplification rules from above:

$$\begin{aligned}
 & (r_1 + \emptyset) \cdot \epsilon + ((\epsilon + r_2) + r_3) \cdot (r_4 \cdot \emptyset) \\
 \equiv & \underbrace{(r_1 + \emptyset) \cdot \epsilon + ((\epsilon + r_2) + r_3) \cdot \emptyset} \\
 \equiv & \underbrace{(r_1 + \emptyset) \cdot \epsilon} + \emptyset \\
 \equiv & \underbrace{(r_1 + \emptyset)} + \emptyset \\
 \equiv & \underbrace{r_1 + \emptyset} \\
 \equiv & r_1
 \end{aligned}$$

In each step, I underlined where a simplification rule is applied. Our matching algorithm in the next section will often generate such “useless”  $\epsilon$ s and  $\emptyset$ s, therefore simplifying them away will make the algorithm quite a bit faster.

## The Matching Algorithm

The algorithm we will define below consists of two parts. One is the function *nullable* which takes a regular expression as argument and decides whether it can match the empty string (this means it returns a boolean in Scala). This can be easily defined recursively as follows:

$$\begin{aligned}
 nullable(\emptyset) & \stackrel{\text{def}}{=} false \\
 nullable(\epsilon) & \stackrel{\text{def}}{=} true \\
 nullable(c) & \stackrel{\text{def}}{=} false \\
 nullable(r_1 + r_2) & \stackrel{\text{def}}{=} nullable(r_1) \vee nullable(r_2) \\
 nullable(r_1 \cdot r_2) & \stackrel{\text{def}}{=} nullable(r_1) \wedge nullable(r_2) \\
 nullable(r^*) & \stackrel{\text{def}}{=} true
 \end{aligned}$$

The idea behind this function is that the following property holds:

$$nullable(r) \text{ if and only if } [] \in L(r)$$

Note on the left-hand side of the if-and-only-if we have a function we can implement; on the right we have its specification (which we cannot implement in a programming language).

The other function of our matching algorithm calculates a *derivative* of a regular expression. This is a function which will take a regular expression, say  $r$ , and a character, say  $c$ , as argument and returns a new regular expression. Be careful that the intuition behind this function is not so easy to grasp on first reading. Essentially this function solves the following problem: if  $r$  can match a string of the form  $c :: s$ , what does the regular expression look like that can match just  $s$ ? The definition of this function is as follows:

$$\begin{aligned}
\text{der } c (\emptyset) &\stackrel{\text{def}}{=} \emptyset \\
\text{der } c (\epsilon) &\stackrel{\text{def}}{=} \emptyset \\
\text{der } c (d) &\stackrel{\text{def}}{=} \text{if } c = d \text{ then } \epsilon \text{ else } \emptyset \\
\text{der } c (r_1 + r_2) &\stackrel{\text{def}}{=} \text{der } c r_1 + \text{der } c r_2 \\
\text{der } c (r_1 \cdot r_2) &\stackrel{\text{def}}{=} \text{if } \text{nullable}(r_1) \\
&\quad \text{then } (\text{der } c r_1) \cdot r_2 + \text{der } c r_2 \\
&\quad \text{else } (\text{der } c r_1) \cdot r_2 \\
\text{der } c (r^*) &\stackrel{\text{def}}{=} (\text{der } c r) \cdot (r^*)
\end{aligned}$$

The first two clauses can be rationalised as follows: recall that *der* should calculate a regular expression so that given the “input” regular expression can match a string of the form  $c :: s$ , we want a regular expression for  $s$ . Since neither  $\emptyset$  nor  $\epsilon$  can match a string of the form  $c :: s$ , we return  $\emptyset$ . In the third case we have to make a case-distinction: In case the regular expression is  $c$ , then clearly it can recognise a string of the form  $c :: s$ , just that  $s$  is the empty string. Therefore we return the  $\epsilon$ -regular expression. In the other case we again return  $\emptyset$  since no string of the  $c :: s$  can be matched. Next come the recursive cases, which are a bit more involved. Fortunately, the  $+$ -case is still relatively straightforward: all strings of the form  $c :: s$  are either matched by the regular expression  $r_1$  or  $r_2$ . So we just have to recursively call *der* with these two regular expressions and compose the results again with  $+$ . Makes sense? The  $\cdot$ -case is more complicated: if  $r_1 \cdot r_2$  matches a string of the form  $c :: s$ , then the first part must be matched by  $r_1$ . Consequently, it makes sense to construct the regular expression for  $s$  by calling *der* with  $r_1$  and “appending”  $r_2$ . There is however one exception to this simple rule: if  $r_1$  can match the empty string, then all of  $c :: s$  is matched by  $r_2$ . So in case  $r_1$  is nullable (that is can match the empty string) we have to allow the choice  $\text{der } c r_2$  for calculating the regular expression that can match  $s$ . Therefore we have to add the regular expression  $\text{der } c r_2$  in the result. The  $*$ -case is again simple: if  $r^*$  matches a string of the form  $c :: s$ , then the first part must be “matched” by a single copy of  $r$ . Therefore we call recursively *der*  $c r$  and “append”  $r^*$  in order to match the rest of  $s$ .

If this did not make sense, here is another way to rationalise the definition of *der* by considering the following operation on sets:

$$\text{Der } c A \stackrel{\text{def}}{=} \{s \mid c :: s \in A\}$$

This operation essentially transforms a set of strings  $A$  by filtering out all strings that do not start with  $c$  and then strips off the  $c$  from all the remaining strings. For example suppose  $A = \{\text{foo}, \text{bar}, \text{frak}\}$  then

$$\text{Der } f A = \{\text{oo}, \text{rak}\} \quad , \quad \text{Der } b A = \{\text{ar}\} \quad \text{and} \quad \text{Der } a A = \emptyset$$

Note that in the last case *Der* is empty, because no string in  $A$  starts with  $a$ . With this operation we can state the following property about *der*:

$$L(\text{der } c r) = \text{Der } c (L(r))$$

This property clarifies what regular expression *der* calculates, namely take the set of strings that *r* can match (that is  $L(r)$ ), filter out all strings not starting with *c* and strip off the *c* from the remaining strings—this is exactly the language that *der c r* can match.

If we want to find out whether the string *abc* is matched by the regular expression  $r_1$  then we can iteratively apply *der* as follows

Input:  $r_1, abc$

- Step 1: build derivative of *a* and  $r_1$       ( $r_2 = der\ a\ r_1$ )
- Step 2: build derivative of *b* and  $r_2$       ( $r_3 = der\ b\ r_2$ )
- Step 3: build derivative of *c* and  $r_3$       ( $r_4 = der\ c\ r_3$ )
- Step 4: the string is exhausted; test      ( $nullable(r_4)$ )  
whether  $r_4$  can recognise the  
empty string

Output: result of this test  $\Rightarrow true$  or *false*

Again the operation *Der* might help to rationalise this algorithm. We want to know whether  $abc \in L(r_1)$ . We do not know yet—but let us assume it is. Then  $Der\ a\ L(r_1)$  builds the set where all the strings not starting with *a* are filtered out. Of the remaining strings, the *a* is stripped off. Then we continue with filtering out all strings not starting with *b* and stripping off the *b* from the remaining strings, that means we build  $Der\ b\ (Der\ a\ (L(r_1)))$ . Finally we filter out all strings not starting with *c* and strip off *c* from the remaining string. This is  $Der\ c\ (Der\ b\ (Der\ a\ (L(r_1))))$ . Now if *abc* was in the original set ( $L(r_1)$ ), then in  $Der\ c\ (Der\ b\ (Der\ a\ (L(r_1))))$  must be the empty string. If not, then *abc* was not in the language we started with.

Our matching algorithm using *der* and *nullable* works similarly, just using regular expression instead of sets. For this we need to extend the notion of derivatives from single characters to strings. This can be done using the following function, taking a string and regular expression as input and a regular expression as output.

$$\begin{aligned} ders\ []\ r &\stackrel{\text{def}}{=} r \\ ders\ (c::s)\ r &\stackrel{\text{def}}{=} ders\ s\ (der\ c\ r) \end{aligned}$$

This function iterates *der* taking one character at the time from the original string until it is exhausted. Having *ders* in place, we can finally define our matching algorithm:

$$matches\ s\ r \stackrel{\text{def}}{=} nullable(ders\ s\ r)$$

and we can claim that

$$matches\ s\ r \quad \text{if and only if} \quad s \in L(r)$$

holds, which means our algorithm satisfies the specification. Of course we can claim many things...whether the claim holds any water is a different question, which for example is the point of the Strand-2 Coursework.

This algorithm was introduced by Janus Brzozowski in 1964. Its main attractions are simplicity and being fast, as well as being easily extendable for other regular expressions such as  $r^{\{n\}}$ ,  $r^?$ ,  $\sim r$  and so on (this is subject of Strand-1 Coursework 1).

## The Matching Algorithm in Scala

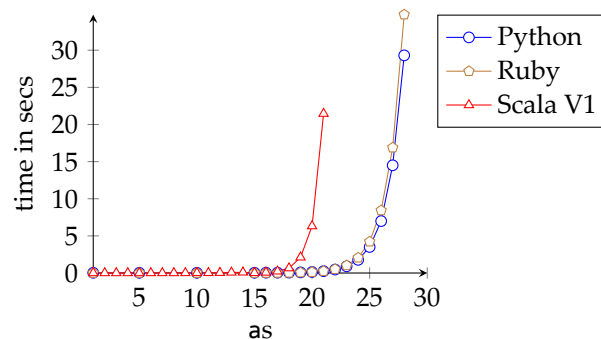
Another attraction of the algorithm is that it can be easily implemented in a functional programming language, like Scala. Given the implementation of regular expressions in Scala shown in the first lecture and handout, the functions and subfunctions for `matches` are shown in Figure 1.

For running the algorithm with our favourite example, the evil regular expression  $a^{\{n\}}a^{\{n\}}$ , we need to implement the optional regular expression and the exactly  $n$ -times regular expression. This can be done with the translations

```
def OPT(r: Rexp) = ALT(r, EMPTY)

def NTIMES(r: Rexp, n: Int) : Rexp = n match {
  case 0 => EMPTY
  case 1 => r
  case n => SEQ(r, NTIMES(r, n - 1))
}
```

Running the matcher with the example, we find it is slightly worse then the matcher in Ruby and Python. Ooops...



Analysing this failure we notice that for  $a^{\{n\}}$  we generate quite big regular expressions:

```

1  def nullable (r: Rexp) : Boolean = r match {
2    case NULL => false
3    case EMPTY => true
4    case CHAR(_) => false
5    case ALT(r1, r2) => nullable(r1) || nullable(r2)
6    case SEQ(r1, r2) => nullable(r1) && nullable(r2)
7    case STAR(_) => true
8  }
9
10 def der (c: Char, r: Rexp) : Rexp = r match {
11   case NULL => NULL
12   case EMPTY => NULL
13   case CHAR(d) => if (c == d) EMPTY else NULL
14   case ALT(r1, r2) => ALT(der(c, r1), der(c, r2))
15   case SEQ(r1, r2) =>
16     if (nullable(r1)) ALT(SEQ(der(c, r1), r2), der(c, r2))
17     else SEQ(der(c, r1), r2)
18   case STAR(r) => SEQ(der(c, r), STAR(r))
19 }
20
21 def ders (s: List[Char], r: Rexp) : Rexp = s match {
22   case Nil => r
23   case c::s => ders(s, der(c, r))
24 }
25
26 def matches(r: Rexp, s: String) : Boolean =
27   nullable(ders(s.toList, r))

```

Figure 1: Scala implementation of the nullable and derivatives functions. These functions are easy to implement in functional languages, because pattern matching and recursion allow us to mimic the mathematical definitions very closely.

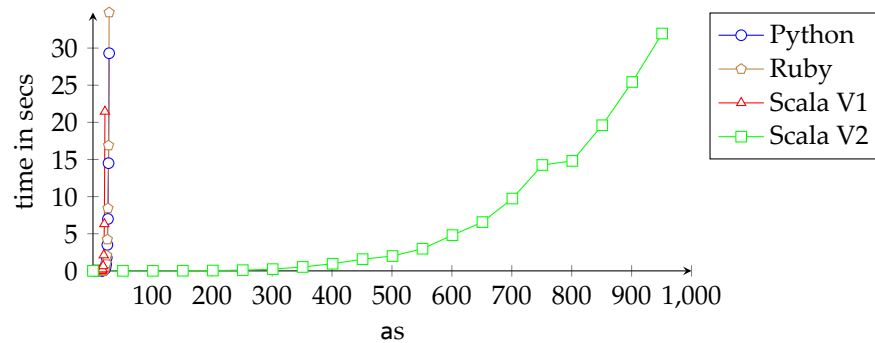
- 1:  $a$
- 2:  $a \cdot a$
- 3:  $a \cdot a \cdot a$
- ...
- 13:  $a \cdot a \cdot a \cdot a \cdot a \cdot a \cdot a \cdot a \cdot a \cdot a \cdot a \cdot a \cdot a$
- ...

Our algorithm traverses such regular expressions at least once every time a derivative is calculated. So having large regular expressions will cause problems. This problem is aggravated by  $a?$  being represented as  $a + \epsilon$ .

We can however fix this by having an explicit constructor for  $r^{\{n\}}$ . In Scala we would introduce a constructor like

```
case class NTIMES(r: Rexp, n: Int) extends Rexp
```

With this fix we have a constant “size” regular expression for our running example no matter how large  $n$  is. This means we have to also add cases for NTIMES in the functions *nullable* and *der*. Does the change have any effect?



Now we are talking business! The modified matcher can within 30 seconds handle regular expressions up to  $n = 950$  before a StackOverflow is raised.

The moral is that our algorithm is rather sensitive to the size of regular expressions it needs to handle. This is of course obvious because both *nullable* and *der* frequently need to traverse the whole regular expression. There seems, however, one more issue for making the algorithm run faster. The derivative function often produces “useless”  $\emptyset$ s and  $\epsilon$ s. To see this, consider  $r = ((a \cdot b) + b)^*$  and the following two derivatives

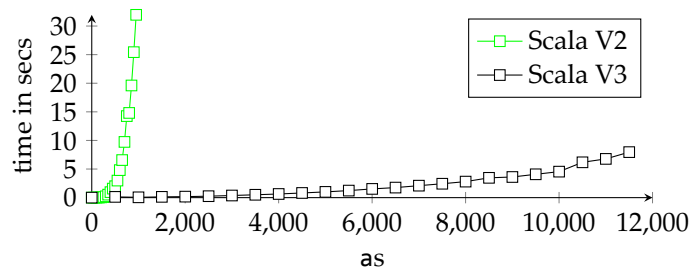
$$\begin{aligned} \text{der } a r &= ((\epsilon \cdot b) + \emptyset) \cdot r \\ \text{der } b r &= ((\emptyset \cdot b) + \epsilon) \cdot r \\ \text{der } c r &= ((\emptyset \cdot b) + \emptyset) \cdot r \end{aligned}$$

If we simplify them according to the simple rules from the beginning, we can replace the right-hand sides by the smaller equivalent regular expressions



$$\begin{aligned} \text{der } a r &\equiv b \cdot r \\ \text{der } b r &\equiv r \\ \text{der } c r &\equiv \emptyset \end{aligned}$$

I leave it to you to contemplate whether such a simplification can have any impact on the correctness of our algorithm (will it change any answers?). Figure 2 gives a simplification function that recursively traverses a regular expression and simplifies it according to the rules given at the beginning. There are only rules for  $+$ ,  $\cdot$  and  $n$ -times (the latter because we added it in the second version of our matcher). There is no rule for a star, because empirical data and also a little thought showed that simplifying under a star is a waste of computation time. The simplification function will be called after every derivation. This additional step removes all the “junk” the derivative function introduced. Does this improve the speed? You bet!!



## Proofs

You might not like doing proofs. But they serve a very important purpose in Computer Science: When can we be sure that our algorithms match their specification. We can try to test the algorithms, but that often overlooks corner cases and also often an exhaustive testing is impossible (since there are often infinitely many inputs). Proofs allow us to ensure that an algorithm meets its specification.

For the programs we look at in this module, the proofs will mostly by some form of induction. Remember that regular expressions are defined as

$r ::= \emptyset$	null
$\epsilon$	empty string / "" / []
$c$	single character
$r_1 + r_2$	alternative / choice
$r_1 \cdot r_2$	sequence
$r^*$	star (zero or more)

If you want to show a property  $P(r)$  for all regular expressions  $r$ , then you have to follow essentially the recipe:

- $P$  has to hold for  $\emptyset$ ,  $\epsilon$  and  $c$  (these are the base cases).

```

1  def simp(r: Rexp): Rexp = r match {
2    case ALT(r1, r2) => {
3      val r1s = simp(r1)
4      val r2s = simp(r2)
5      (r1s, r2s) match {
6        case (NULL, _) => r2s
7        case (_, NULL) => r1s
8        case _ => if (r1s == r2s) r1s else ALT(r1s, r2s)
9      }
10   }
11  case SEQ(r1, r2) => {
12    val r1s = simp(r1)
13    val r2s = simp(r2)
14    (r1s, r2s) match {
15      case (NULL, _) => NULL
16      case (_, NULL) => NULL
17      case (EMPTY, _) => r2s
18      case (_, EMPTY) => r1s
19      case _ => SEQ(r1s, r2s)
20    }
21  }
22  case NTIMES(r, n) => NTIMES(simp(r), n)
23  case r => r
24 }
25
26 def ders (s: List[Char], r: Rexp) : Rexp = s match {
27   case Nil => r
28   case c::s => ders(s, simp(der(c, r)))
29 }

```

Figure 2: The simplification function and modified ders-function; this function now calls der first, but then simplifies the resulting derivative regular expressions, see Line 28.

- $P$  has to hold for  $r_1 + r_2$  under the assumption that  $P$  already holds for  $r_1$  and  $r_2$ .
- $P$  has to hold for  $r_1 \cdot r_2$  under the assumption that  $P$  already holds for  $r_1$  and  $r_2$ .
- $P$  has to hold for  $r^*$  under the assumption that  $P$  already holds for  $r$ .

A simple proof is for example showing the following property:

$$\text{nullable}(r) \text{ if and only if } \epsilon \in L(r)$$

Let us say that this property is  $P(r)$ , then the first case we need to check is whether  $P(\emptyset)$  (see recipe above). So we have to show that

$$\text{nullable}(\emptyset) \text{ if and only if } \epsilon \in L(\emptyset)$$

whereby  $\text{nullable}(\emptyset)$  is by definition of the function *nullable* always *false*. We also have that  $L(\emptyset)$  is by definition  $\{\}$ . It is impossible that the empty string  $\epsilon$  is in the empty set. Therefore also the right-hand side is false. Consequently we verified this case. We would still need to do this for  $P(\epsilon)$  and  $P(c)$ . I leave this to you.