

Handout 4 (Sulzmann & Lu Algorithm)

So far our algorithm based on derivatives was only able to say yes or no depending on whether a string was matched by regular expression or not. Often a more interesting question is to find out *how* a regular expression matched a string? Answering this question will also help us with the problem we are after, namely tokenising an input string.

The algorithm we will be looking at was designed by Sulzmann & Lu in a rather recent paper (from 2014). A link to it is provided on KEATS, in case you are interested.¹ In order to give an answer for *how* a regular expression matches a string, Sulzmann and Lu introduce *values*. A value will be the output of the algorithm whenever the regular expression matches the string. If the string does not match the string, an error will be raised. Since the first phase of the algorithm by Sulzmann & Lu is identical to the derivative based matcher from the first coursework, the function *nullable* will be used to decide whether a string is matched by a regular expression. If *nullable* says yes, then values are constructed that reflect how the regular expression matched the string.

The definitions for values is given below. They are shown together with the regular expressions r to which they correspond:

regular expressions	values
$r ::= \emptyset$	$v ::=$
ϵ	<i>Empty</i>
c	<i>Char(c)</i>
$r_1 \cdot r_2$	<i>Seq(v₁, v₂)</i>
$r_1 + r_2$	<i>Left(v)</i>
r^*	<i>Right(v)</i>
	[<i>v₁... v_n</i>]

The reason is that there is a very strong correspondence between them. There is no value for the \emptyset regular expression, since it does not match any string. Otherwise there is exactly one value corresponding to each regular expression with the exception of $r_1 + r_2$ where there are two values, namely *Left(v)* and *Right(v)* corresponding to the two alternatives. Note that r^* is associated with a list of values, one for each copy of r that was needed to match the string. This means we might also return the empty list [], if no copy was needed in case of r^* . For sequence, there is exactly one value, composed of two component values (v_1 and v_2).

To emphasise the connection between regular expressions and values, I have in my implementation the convention that regular expressions and values have the same name, except that regular expressions are written entirely with upper-case letters, while values just start with a single upper-case character and the rest are lower-case letters. My definition of regular expressions and values in

¹In my humble opinion this is an interesting instance of the research literature: it contains a very neat idea, but its presentation is rather sloppy. In earlier versions of their paper, a King's student and I found several rather annoying typos in their examples and definitions.

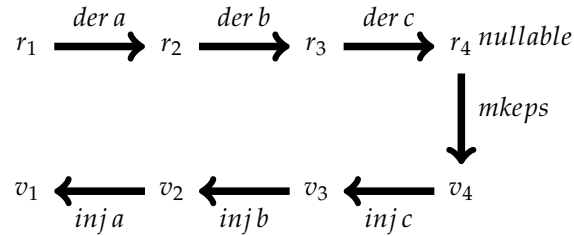


Figure 1: The two phases of the algorithm by Sulzmann & Lu.

Scala is shown below. I use this in the REPL of Scala; when I use the Scala compiler I need to rename some constructors, because Scala on Macs does not like classes that are called `EMPTY` and `Empty`.

```

abstract class Rexp
case object NULL extends Rexp
case object EMPTY extends Rexp
case class CHAR(c: Char) extends Rexp
case class ALT(r1: Rexp, r2: Rexp) extends Rexp
case class SEQ(r1: Rexp, r2: Rexp) extends Rexp
case class STAR(r: Rexp) extends Rexp

abstract class Val
case object Empty extends Val
case class Char(c: Char) extends Val
case class Seq(v1: Val, v2: Val) extends Val
case class Left(v: Val) extends Val
case class Right(v: Val) extends Val
case class Stars(vs: List[Val]) extends Val
  
```

Graphically the algorithm by Sulzmann & Lu can be illustrated by the picture in Figure 1 where the path from the left to the right involving *der/nullable* is the first phase of the algorithm and *mkeps/inj*, the path from right to left, the second phase. This picture shows the steps required when a regular expression, say r_1 , matches the string abc . We first build the three derivatives (according to a , b and c). We then use *nullable* to find out whether the resulting regular expression can match the empty string. If yes, we call the function *mkeps*.

The *mkeps* function calculates a value for how a regular expression has matched the empty string. Its definition is as follows:

$$\begin{aligned}
mkeps(\epsilon) &\stackrel{\text{def}}{=} \text{Empty} \\
mkeps(r_1 + r_2) &\stackrel{\text{def}}{=} \text{if } \text{nullable}(r_1) \\
&\quad \text{then } \text{Left}(mkeps(r_1)) \\
&\quad \text{else } \text{Right}(mkeps(r_2)) \\
mkeps(r_1 \cdot r_2) &\stackrel{\text{def}}{=} \text{Seq}(mkeps(r_1), mkeps(r_2)) \\
mkeps(r^*) &\stackrel{\text{def}}{=} []
\end{aligned}$$

There are no cases for \emptyset and c , since these regular expression cannot match the empty string. Note also that in case of alternatives we give preference to the regular expression on the left-hand side. This will become important later on.

The second phase of the algorithm is organised so that it will calculate a value for how the derivative regular expression has matched a string. For this we need a function that reverses this “chopping off” for values which we did in the first phase for derivatives. The corresponding function is called *inj* for injection. This function takes three arguments: the first one is a regular expression for which we want to calculate the value, the second is the character we want to inject and the third argument is the value where we will inject the character into. The result of this function is a new value. The definition of *inj* is as follows:

$$\begin{aligned}
inj(c) c \text{ Empty} &\stackrel{\text{def}}{=} \text{Char } c \\
inj(r_1 + r_2) c \text{ Left}(v) &\stackrel{\text{def}}{=} \text{Left}(inj r_1 c v) \\
inj(r_1 + r_2) c \text{ Right}(v) &\stackrel{\text{def}}{=} \text{Right}(inj r_2 c v) \\
inj(r_1 \cdot r_2) c \text{ Seq}(v_1, v_2) &\stackrel{\text{def}}{=} \text{Seq}(inj r_1 c v_1, v_2) \\
inj(r_1 \cdot r_2) c \text{ Left}(\text{Seq}(v_1, v_2)) &\stackrel{\text{def}}{=} \text{Seq}(inj r_1 c v_1, v_2) \\
inj(r_1 \cdot r_2) c \text{ Right}(v) &\stackrel{\text{def}}{=} \text{Seq}(mkeps(r_1), inj r_2 c v) \\
inj(r^*) c \text{ Seq}(v, vs) &\stackrel{\text{def}}{=} inj r c v :: vs
\end{aligned}$$

This definition is by recursion on the regular expression and by analysing the shape of the values. Therefore there are, for example, three cases for sequence regular expressions (for all possible shapes of the value). The last clause for the star regular expression returns a list where the first element is *inj r c v* and the other elements are *vs*. That means *_ :: _* should be read as list cons.

To understand what is going on, it might be best to do some example calculations and compare them with Figure 1. For this note that we have not yet dealt with the need of simplifying regular expressions (this will be a topic on its own later). Suppose the regular expression is $a \cdot (b \cdot c)$ and the input string is *abc*. The derivatives from the first phase are as follows:

$$\begin{aligned}
r_1: & a \cdot (b \cdot c) \\
r_2: & \epsilon \cdot (b \cdot c) \\
r_3: & (\emptyset \cdot (b \cdot c)) + (\epsilon \cdot c) \\
r_4: & (\emptyset \cdot (b \cdot c)) + ((\emptyset \cdot c) + \epsilon)
\end{aligned}$$

According to the simple algorithm, we would test whether r_4 is nullable, which in this case it indeed is. This means we can use the function *mkeps* to calculate a value for how r_4 was able to match the empty string. Remember that this function gives preference for alternatives on the left-hand side. However there is only ϵ on the very right-hand side of r_4 that matches the empty string. Therefore *mkeps* returns the value

$$v_4 : \text{Right}(\text{Right}(\text{Empty}))$$

If there had been a ϵ on the left, then *mkeps* would have returned something of the form *Left*(...). The point is that from this value we can directly read off which part of r_4 matched the empty string: take the right-alternative first, and then the right-alternative again.

Next we have to “inject” the last character, that is c in the running example, into this value v_4 in order to calculate how r_3 could have matched the string c . According to the definition of *inj* we obtain

$$v_3 : \text{Right}(\text{Seq}(\text{Empty}, \text{Char}(c)))$$

This is the correct result, because r_3 needs to use the right-hand alternative, and then ϵ needs to match the empty string and c needs to match c . Next we need to inject back the letter b into v_3 . This gives

$$v_2 : \text{Seq}(\text{Empty}, \text{Seq}(\text{Char}(b), \text{Char}(c)))$$

which is again the correct result for how r_2 matched the string bc . Finally we need to inject back the letter a into v_2 giving the final result

$$v_1 : \text{Seq}(\text{Char}(a), \text{Seq}(\text{Char}(b), \text{Char}(c)))$$

This now corresponds to how the regular expression $a \cdot (b \cdot c)$ matched the string abc .

There are a few auxiliary functions that are of interest when analysing this algorithm. One is called *flatten*, written $|_|_$, which extracts the string “underlying” a value. It is defined recursively as

$$\begin{aligned} |\text{Empty}| &\stackrel{\text{def}}{=} [] \\ |\text{Char}(c)| &\stackrel{\text{def}}{=} [c] \\ |\text{Left}(v)| &\stackrel{\text{def}}{=} |v| \\ |\text{Right}(v)| &\stackrel{\text{def}}{=} |v| \\ |\text{Seq}(v_1, v_2)| &\stackrel{\text{def}}{=} |v_1| @ |v_2| \\ |[v_1, \dots, v_n]| &\stackrel{\text{def}}{=} |v_1| @ \dots @ |v_n| \end{aligned}$$

Using *flatten* we can see what are the strings behind the values calculated by *mkeps* and *inj*. In our running example:

$$\begin{array}{l}
|v_4|: \quad \square \\
|v_3|: \quad c \\
|v_2|: \quad bc \\
|v_1|: \quad abc
\end{array}$$

This indicates that *inj* indeed is injecting, or adding, back a character into the value. If we continue until all characters are injected back, we have a value that can indeed say how the string *abc* was matched.

There is a problem, however, with the described algorithm so far: it is very slow. We need to include the simplification from Lecture 2. This is what we shall do next.

Simplification

Generally the matching algorithms based on derivatives do poorly unless the regular expressions are simplified after each derivative step. But this is a bit more involved in the algorithm of Sulzmann & Lu. So what follows might require you to read several times before it makes sense and also might require that you do some example calculations yourself. As a first example consider the last derivation step in our earlier example:

$$r_4 = \text{der } c r_3 = (\emptyset \cdot (b \cdot c)) + ((\emptyset \cdot c) + \epsilon)$$

Simplifying this regular expression would just give us ϵ . Running *mkeys* with this regular expression as input, however, would then provide us with *Empty* instead of *Right(Right(Empty))* that was obtained without the simplification. The problem is we need to recreate this more complicated value, rather than just return *Empty*.

This will require what I call *rectification functions*. They need to be calculated whenever a regular expression gets simplified. Rectification functions take a value as argument and return a (rectified) value. Let us first take a look again at our simplification rules:

$$\begin{array}{l}
r \cdot \emptyset \mapsto \emptyset \\
\emptyset \cdot r \mapsto \emptyset \\
r \cdot \epsilon \mapsto r \\
\epsilon \cdot r \mapsto r \\
r + \emptyset \mapsto r \\
\emptyset + r \mapsto r \\
r + r \mapsto r
\end{array}$$

Applying them to r_4 will require several nested simplifications in order end up with just ϵ . However, it is possible to apply them in a depth-first, or inside-out, manner in order to calculate this simplified regular expression.

The rectification we can implement by letting *simp* return not just a (simplified) regular expression, but also a rectification function. Let us consider the

alternative case, $r_1 + r_2$, first. By going depth-first, we first simplify the component regular expressions r_1 and r_2 . This will return simplified versions (if they can be simplified), say r_{1s} and r_{2s} , but also two rectification functions f_{1s} and f_{2s} . We need to assemble them in order to obtain a rectified value for $r_1 + r_2$. In case r_{1s} simplified to \emptyset , we continue the derivative calculation with r_{2s} . The Sulzmann & Lu algorithm would return a corresponding value, say v_{2s} . But now this value needs to be “rectified” to the value

$$Right(v_{2s})$$

The reason is that we look for the value that tells us how $r_1 + r_2$ could have matched the string, not just r_2 or r_{2s} . Unfortunately, this is still not the right value in general because there might be some simplifications that happened inside r_2 and for which the simplification function returned also a rectification function f_{2s} . So in fact we need to apply this one too which gives

$$Right(f_{2s}(v_{2s}))$$

This is now the correct, or rectified, value. Since the simplification will be done in the first phase of the algorithm, but the rectification needs to be done to the values in the second phase, it is advantageous to calculate the rectification as a function, remember this function and then apply the value to this function during the second phase. So if we want to implement the rectification as function, we would need to return

$$\lambda v. Right(f_{2s}(v))$$

which is the lambda-calculus notation for a function that expects a value v and returns everything after the dot where v is replaced by whatever value is given.

Let us package this idea with rectification functions into a single function (still only considering the alternative case):

```

simp(r):
  case r = r1 + r2
    let (r1s, f1s) = simp(r1)
        (r2s, f2s) = simp(r2)
    case r1s = ∅: return (r2s, λv. Right(f2s(v)))
    case r2s = ∅: return (r1s, λv. Left(f1s(v)))
    case r1s = r2s: return (r1s, λv. Left(f1s(v)))
    otherwise: return (r1s + r2s, falt(f1s, f2s))

```

We first recursively call the simplification with r_1 and r_2 . This gives simplified regular expressions, r_{1s} and r_{2s} , as well as two rectification functions f_{1s} and f_{2s} . We next need to test whether the simplified regular expressions are \emptyset so as to make further simplifications. In case r_{1s} is \emptyset , then we can return r_{2s} (the other alternative). However for this we need to build a corresponding rectification function, which as said above is

$$\lambda v. \text{Right}(f_{2s}(v))$$

The case where $r_{2s} = \emptyset$ is similar: We return r_{1s} and rectify with $\text{Left}(_)$ and the other calculated rectification function f_{1s} . This gives

$$\lambda v. \text{Left}(f_{1s}(v))$$

The next case where $r_{1s} = r_{2s}$ can be treated like the one where $r_{2s} = \emptyset$. We return r_{1s} and rectify with $\text{Left}(_)$ and so on.

The otherwise-case is slightly more complicated. In this case neither r_{1s} nor r_{2s} are \emptyset and also $r_{1s} \neq r_{2s}$, which means no further simplification can be applied. Accordingly, we return $r_{1s} + r_{2s}$ as the simplified regular expression. In principle we also do not have to do any rectification, because no simplification was done in this case. But this is actually not true: There might have been simplifications inside r_{1s} and r_{2s} . We therefore need to take into account the calculated rectification functions f_{1s} and f_{2s} . We can do this by defining a rectification function f_{alt} which takes two rectification functions as arguments and applies them according to whether the value is of the form $\text{Left}(_)$ or $\text{Right}(_)$:

$$f_{alt}(f_1, f_2) \stackrel{\text{def}}{=} \begin{array}{l} \lambda v. \text{case } v = \text{Left}(v'): \text{return } \text{Left}(f_1(v')) \\ \text{case } v = \text{Right}(v'): \text{return } \text{Right}(f_2(v')) \end{array}$$

The other interesting case with simplification is the sequence case. In this case the main simplification function is as follows

$$\begin{array}{l} \text{simp}(r): \quad \text{(continued)} \\ \text{case } r = r_1 \cdot r_2 \\ \quad \text{let } (r_{1s}, f_{1s}) = \text{simp}(r_1) \\ \quad \quad (r_{2s}, f_{2s}) = \text{simp}(r_2) \\ \quad \text{case } r_{1s} = \emptyset: \text{return } (\emptyset, f_{error}) \\ \quad \text{case } r_{2s} = \emptyset: \text{return } (\emptyset, f_{error}) \\ \quad \text{case } r_{1s} = \epsilon: \text{return } (r_{2s}, \lambda v. \text{Seq}(f_{1s}(\text{Empty}), f_{2s}(v))) \\ \quad \text{case } r_{2s} = \epsilon: \text{return } (r_{1s}, \lambda v. \text{Seq}(f_{1s}(v), f_{2s}(\text{Empty}))) \\ \quad \text{otherwise: return } (r_{1s} \cdot r_{2s}, f_{seq}(f_{1s}, f_{2s})) \end{array}$$

whereby in the last line f_{seq} is again pushing the two rectification functions into the two components of the Seq-value:

$$f_{seq}(f_1, f_2) \stackrel{\text{def}}{=} \lambda v. \text{case } v = \text{Seq}(v_1, v_2): \text{return } \text{Seq}(f_1(v_1), f_2(v_2))$$

Note that in the case of $r_{1s} = \emptyset$ (similarly r_{2s}) we use the function f_{error} for rectification. If you think carefully, then you will realise that this function will actually never be called. This is because a sequence with \emptyset will never recognise any string and therefore the second phase of the algorithm would never be called. The simplification function still expects us to give a function. So in

```

simp(r):
  case r = r1 + r2
    let (r1s, f1s) = simp(r1)
        (r2s, f2s) = simp(r2)
    case r1s = ∅: return (r2s, λv. Right(f2s(v)))
    case r2s = ∅: return (r1s, λv. Left(f1s(v)))
    case r1s = r2s: return (r1s, λv. Left(f1s(v)))
    otherwise: return (r1s + r2s, falt(f1s, f2s))

  case r = r1 · r2
    let (r1s, f1s) = simp(r1)
        (r2s, f2s) = simp(r2)
    case r1s = ∅: return (∅, ferror)
    case r2s = ∅: return (∅, ferror)
    case r1s = ε: return (r2s, λv. Seq(f1s(Empty), f2s(v)))
    case r2s = ε: return (r1s, λv. Seq(f1s(v), f2s(Empty)))
    otherwise: return (r1s · r2s, fseq(f1s, f2s))

  otherwise:
    return (r, λv. v)

```

Figure 2: The simplification function that returns a simplified regular expression and a rectification function.

my own implementation I just returned a function that raises an error. In the case where $r_{1s} = \epsilon$ (similarly r_{2s}) we have to create a sequence where the first component is a rectified version of *Empty*. Therefore we call f_{1s} with *Empty*.

Since we only simplify regular expressions of the form $r_1 + r_2$ and $r_1 \cdot r_2$ we do not have to do anything else in the remaining cases. The rectification function will be just the identity, which in lambda-calculus terms is just

$$\lambda v. v$$

This completes the high-level version of the simplification function, which is also shown again in Figure 2. This can now be used in a *lexing function* as follows:

$$\begin{aligned}
\text{lex } r \ [] &\stackrel{\text{def}}{=} \text{if } \text{nullable}(r) \text{ then } \text{mkeys}(r) \\
&\quad \text{else } \text{error} \\
\text{lex } r \ c::s &\stackrel{\text{def}}{=} \text{let } (r_{\text{simp}}, f_{\text{rect}}) = \text{simp}(\text{der}(c, r)) \\
&\quad \text{in } \text{j } r \ c \ f_{\text{rect}}(\text{lex } r_{\text{simp}} \ s)
\end{aligned}$$

This corresponds to the *matches* function we have seen in earlier lectures. In the first clause we are given an empty string, [], and need to test whether the regular expression is *nullable*. If yes, we can proceed normally and just return the value calculated by *mkeys*. The second clause is for strings where the first character

is c , say, and the rest of the string is s . We first build the derivative of r with respect to c ; simplify the resulting regular expression. We continue lexing with the simplified regular expression and the string s . Whatever will be returned as value, we still need to rectify using the f_{rect} from the simplification and finally inject c back into the (rectified) value.

Records

Remember we wanted to tokenize input strings, that means splitting strings into their “word” components. Furthermore we want to classify each token as being a keyword or identifier and so on. For this one more feature will be required, which I call a *record* regular expression. While values encode how a regular expression matches a string, records can be used to “focus” on some particular parts of the regular expression and “forget” about others.

Let us look at an example. Suppose you have the regular expression $a \cdot b + a \cdot c$. Clearly this regular expression can only recognise two strings. But suppose you are not interested whether it can recognise ab or ac , but rather if it matched, then what was the last character of the matched string...either b or c . You can do this by annotating the regular expression with a record, written in general $(x : r)$, where x is just an identifier (in my implementation a plain string) and r is a regular expression. A record will be regarded as a regular expression. The extended definition in Scala therefore looks as follows:

```

1 abstract class Rexp
2 case object NULL extends Rexp
3 case object EMPTY extends Rexp
4 case class CHAR(c: Char) extends Rexp
5 case class ALT(r1: Rexp, r2: Rexp) extends Rexp
6 case class SEQ(r1: Rexp, r2: Rexp) extends Rexp
7 case class STAR(r: Rexp) extends Rexp
8 case class REC(x: String, r: Rexp) extends Rexp

```

Since we regard records as regular expressions we need to extend the functions *nullable* and *der*. Similarly *mkeys* and *inj* need to be extended. This means we also need to extend the definition of values, which in Scala looks as follows:

```

1 abstract class Val
2 case object Empty extends Val
3 case class Char(c: Char) extends Val
4 case class Seq(v1: Val, v2: Val) extends Val
5 case class Left(v: Val) extends Val
6 case class Right(v: Val) extends Val
7 case class Stars(vs: List[Val]) extends Val
8 case class Rec(x: String, v: Val) extends Val

```

Let us now look at the purpose of records more closely and let us return to our question whether the string terminated in a b or c . We can do this as follows: we annotate the regular expression $ab + ac$ with a record as follows

$$a \cdot (x : b) + a \cdot (x : c)$$

This regular expression can still only recognise the strings ab and ac , but we can now use a function that takes a value and returns all records. I call this function env for environment...it builds a list of identifiers associated with a string. This function can be defined as follows:

$$\begin{aligned}
env(Empty) &\stackrel{\text{def}}{=} [] \\
env(Char(c)) &\stackrel{\text{def}}{=} [] \\
env(Left(v)) &\stackrel{\text{def}}{=} env(v) \\
env(Right(v)) &\stackrel{\text{def}}{=} env(v) \\
env(Seq(v_1, v_2)) &\stackrel{\text{def}}{=} env(v_1) @ env(v_2) \\
env([v_1, \dots, v_n]) &\stackrel{\text{def}}{=} env(v_1) @ \dots @ env(v_n) \\
env(Rec(x : v)) &\stackrel{\text{def}}{=} (x : |v|) :: env(v)
\end{aligned}$$

where in the last clause we use the $flatten$ function defined earlier. As can be seen, the function env “picks” out all underlying strings where a record is given. Since there can be more than one, the environment will potentially contain many “records”. If we now postprocess the value calculated by lex extracting all records using env , we can answer the question whether the last element in the string was an b or a c . Lets see this in action: if we use $a \cdot b + a \cdot c$ and ac the calculated value will be

$$Right(Seq(Char(a), Char(c)))$$

If we use instead $a \cdot (x : b) + a \cdot (x : c)$ and use the env function to extract the recording for x we obtain

$$[(x : c)]$$

If we had given the string ab instead, then the record would have been $[(x : b)]$. The fun starts if we iterate this. Consider the regular expression

$$(a \cdot (x : b) + a \cdot (y : c))^*$$

and the string $ababacabacab$. This string is clearly matched by the regular expression, but we are only interested in the sequence of bs and cs . Using env we obtain

$$[(x : b), (x : b), (y : c), (x : b), (y : c), (x : b)]$$

While this feature might look silly, it is in fact quite useful. For example if we want to match the name of an email we might use the regular expression

$$(name : [a-z0-9_.-]^+) \cdot @ \cdot (domain : [a-z0-9.-]^+) \cdot \dots \cdot (top_level : [a-z.]{2,6})$$

Then if we match the email address

`christian.urban@kcl.ac.uk`

we can use the *env* function and find out what the name, domain and top-level part of the email address are:

```
[(name : christian.urban), (domain : kcl), (top_level : ac.uk)]
```

As you will see in the next lecture, this is now all we need to tokenise an input string and classify each token.