# Regexp Module

smlhelp.github.io – Auxiliary Library

Jacob Neumann
July 2021

# Contents

**Note 1**

Throughout, `Sigma` is assumed to be some equality type.

**Note 2**

This document assumes the definitions and declarations of https://github.com/smlhelp/aux-library/blob/main/documentation/language.pdf, and assumes that the `Language` structure has been opened.

In particular, we have

**Defn. 3**

aux-library: Language.sml

```
31    type 'S language = 'S list -> bool
```

A total function `D : Sigma list -> bool` is said to **compute** a set $L$ of values of type `Sigma list` if

$$D \ cs \Longrightarrow true \qquad \text{iff} \qquad cs \in L$$

i.e. we can run `D` on any value `cs : Sigma list` and get back either `true` or `false`, which correctly indicates that either $cs \in L$ or $cs \notin L$, respectively.

# 1    Preliminary Declarations & Main Definitions

## Declaration 4

aux-library: Regexp.sml

```
29    datatype ''S regexp =
30        Zero
31        | One
32        | Const of ''S
33        | Plus of ''S regexp * ''S regexp
34        | Times of ''S regexp * ''S regexp
35        | Star of ''S regexp
```

## Defn. 5 (Language)

Given a value `R : Sigma regexp`, define the **language of** `R` to be the set of values $\mathcal{L}(\texttt{R}) \subseteq \textsf{Values}(\texttt{Sigma list})$ given recursively by:

$$
\begin{aligned}
\mathcal{L}(\texttt{Zero}) &= \emptyset \\
\mathcal{L}(\texttt{One}) &= \{\texttt{[]}\} \\
\mathcal{L}(\texttt{Const(c)}) &= \{\texttt{[c]}\} \\
\mathcal{L}(\texttt{Plus(r1,r2)}) &= \mathcal{L}(\texttt{r1}) \cup \mathcal{L}(\texttt{r2}) \\
\mathcal{L}(\texttt{Times(r1,r2)}) &= \{\texttt{v}_1 \texttt{@v}_2 \mid \texttt{v}_1 \in \mathcal{L}(\texttt{r1}) \text{ and } \texttt{v}_2 \in \mathcal{L}(\texttt{r2})\} \\
\mathcal{L}(\texttt{Star(r)}) &= \{\texttt{v}_1 \texttt{@v}_2 \texttt{@ ... @v}_n \mid n \in \mathbb{N},\ \texttt{v}_1, \texttt{v}_2, \ldots, \texttt{v}_n \in \mathcal{L}(\texttt{r})\}
\end{aligned}
$$

## Defn. 6

A total function `D : Sigma regexp -> Sigma language` is said to **compute the regular expression decision problem** (over the alphabet `Sigma`) if

$$\texttt{D R cs} \implies \texttt{true} \quad \text{if and only if} \quad \texttt{cs} \in \mathcal{L}(\texttt{R})$$

for all `R : Sigma regexp` and all `cs : Sigma list`.

## Defn. 7

A pair `(p,s) : Sigma list * Sigma list` is said to be a **splitting** of `cs : Sigma list` if

$$\texttt{cs} \cong \texttt{p @ s}.$$

In such a splitting, `p` is called the **prefix** and `s` the **suffix**.

## Declaration 8

aux-library: Regexp.sml

```
47    exception NoMatch
```

### Defn. 9

For any type `t`, we'll say that a function

$$k \ : \ \text{Sigma list} \ * \ \text{Sigma list} \ \text{->} \ t$$

is **almost total** if, for all `(p,s) : Sigma list * Sigma list`, either

$$k(p,s) \hookrightarrow v \text{ for some value } v \qquad \text{or} \qquad k(p,s) \text{ raises NoMatch.}$$

If `k(p,s)` $\hookrightarrow$ `v`, then we say that **k accepts (p,s) with result v**.

# 2   Matching Specification

## Declaration 10 (Depth)

Define the **depth** of `R : Sigma regexp` recursively by

aux-library: Regexp.sml

```
37    fun depth Zero = 0
38      | depth One = 0
39      | depth (Const(_)) = 0
40      | depth (Plus(R1,R2)) =
41          1 + Int.max(depth R1,depth R2)
42      | depth (Times(R1,R2)) =
43          1 + Int.max(depth R1,depth R2)
44      | depth (Star(R)) =
45          1 + depth R
```

## Value Spec

```
match : ''a regexp
          -> ''a list
          -> (''a list * ''a list -> 'b)
          -> 'b
```

**REQUIRES:** `k` is almost total (Defn. 9)

**ENSURES:**

$$
\texttt{match r cs k} \cong \begin{cases} \texttt{v} & \text{where } (\texttt{p},\texttt{s}) \text{ is a splitting of } \texttt{cs} \text{ such} \\ & \text{that } \texttt{p} \in \mathcal{L}(\texttt{r}) \text{ and } \texttt{k} \text{ accepts } (\texttt{p},\texttt{s}) \\ & \text{with result } \texttt{v}. \\ \texttt{raise NoMatch} & \text{if there is no such } (\texttt{p},\texttt{s}) \end{cases}
$$

## Value Spec

```
LL : ''S regexp -> ''S language
```

**ENSURES:** LL computes the regular expression decision problem: for all `R : Sigma regexp`, `(LL R) : Sigma list -> bool` is a total function such that

$$
\texttt{LL R cs} \Longrightarrow \texttt{true} \quad \text{iff} \quad \texttt{cs} \in \mathcal{L}(\texttt{R})
$$

# 3   Matching Implementation

aux-library: Regexp.sml

```
29   datatype ''S regexp =
30       Zero
31       | One
32       | Const of ''S
33       | Plus of ''S regexp * ''S regexp
34       | Times of ''S regexp * ''S regexp
35       | Star of ''S regexp
36
37   fun depth Zero = 0
38     | depth One = 0
39     | depth (Const(_)) = 0
40     | depth (Plus(R1,R2)) =
41         1 + Int.max(depth R1,depth R2)
42     | depth (Times(R1,R2)) =
43         1 + Int.max(depth R1,depth R2)
44     | depth (Star(R)) =
45         1 + depth R
46
47   exception NoMatch
```

aux-library: Regexp.sml

```sml
49   fun match Zero _ _ = raise NoMatch
50     | match One cs k = k([],cs)
51     | match (Const(c)) [] k = raise NoMatch
52     | match (Const(c)) (c'::cs') k =
53         if c=c'
54         then k([c'], cs')
55         else raise NoMatch
56     | match (Plus(R1,R2)) cs k =
57         (match R1 cs k
58           handle NoMatch => match R2 cs k)
59     | match (Times(R1,R2)) cs k =
60         match R1 cs (fn (res',cs') =>
61           match R2 cs' (fn (res'',cs'') =>
62             k (res'@res'',cs'')))
63     | match (Star(R)) cs k =
64         k([],cs)
65           handle NoMatch =>
66             match R cs (fn (res',cs') =>
67               if (cs = cs')
68               then raise NoMatch
69               else
70                 match (Star(R)) cs' (fn (res'',cs'') =>
71                   k(res'@res'',cs'')))
72
73   val LL = fn R => fn s =>
74     match R s (fn (_,[]) => true | _ => raise NoMatch)
75     handle NoMatch => false
```

# 4   Zero, One, Const Correctness

---

**Lemma 11**

For any almost total `k` and any `cs : Sigma list`,

$$
\texttt{match Zero cs k} \cong
\begin{cases}
\texttt{v} & \text{where } (\texttt{p},\texttt{s}) \text{ is a splitting of } \texttt{cs} \text{ such} \\
& \text{that } \texttt{p} \in \mathcal{L}(\texttt{Zero}) \text{ and } \texttt{k} \text{ accepts } (\texttt{p} \\
& ,\texttt{s}) \text{ with result } \texttt{v}. \\
\texttt{raise NoMatch} & \text{if there is no such } (\texttt{p},\texttt{s})
\end{cases}
$$

---

*Proof.* Since $\mathcal{L}(\texttt{Zero}) = \emptyset$, there cannot be any splitting of any `cs` whose prefix is in $\mathcal{L}(\texttt{Zero})$. Therefore, this spec says that `match Zero cs k` always ought to raise `NoMatch`, which it does.   **Spec Satisfied!**

---

**Cor. 12**

For any almost total `k`, (`fn cs => match Zero cs k`) is almost total.

---

**Lemma 13**

For any almost total `k` and any `cs : Sigma list`,

$$
\texttt{match One cs k} \cong
\begin{cases}
\texttt{v} & \text{where } (\texttt{p},\texttt{s}) \text{ is a splitting of } \texttt{cs} \text{ such} \\
& \text{that } \texttt{p} \in \mathcal{L}(\texttt{One}) \text{ and } \texttt{k} \text{ accepts } (\texttt{p}, \\
& \texttt{s}) \text{ with result } \texttt{v}. \\
\texttt{raise NoMatch} & \text{if there is no such } (\texttt{p},\texttt{s})
\end{cases}
$$

---

*Proof.* Since $\mathcal{L}(\texttt{One}) = \{\texttt{[]}\}$, the only splitting of `cs` which we need to consider is (`[]`, `cs`). But notice that

$$
\texttt{k([],cs)} \cong
\begin{cases}
\texttt{v} & \text{if } \texttt{k} \text{ accepts } (\texttt{[]},\texttt{cs}) \text{ with result } \texttt{v} \\
\texttt{raise NoMatch} & \text{if } \texttt{k} \text{ rejects } (\texttt{[]},\texttt{cs})
\end{cases}
$$

so we have   **Spec Satisfied!**   by putting `match One cs k = k([],cs)`.            □

---

**Cor. 14**

For any almost total `k`, (`fn cs => match One cs k`) is almost total.

**Lemma 15**

For any almost total `k`, any `cs : Sigma list`, and any value `c : Sigma`

$$
\texttt{match (Const c) cs k} \cong
\begin{cases}
\texttt{v} & \text{where } \texttt{(p,s)} \text{ is a split-} \\
& \text{ting of } \texttt{cs} \text{ such that } \texttt{p} \in \\
& \mathcal{L}(\texttt{Const c}) \text{ and k accepts (} \\
& \texttt{p,s}) \text{ with result v.} \\
\texttt{raise NoMatch} & \text{if there is no such } \texttt{(p,s)}
\end{cases}
$$

*Proof.* Begin by recalling that $\mathcal{L}(\texttt{Const c}) = \{\texttt{[c]}\}$. Break into two cases: `cs=[]` and `cs=c'::cs'`.

If `cs=[]`, then there is no prefix of `cs` in $\mathcal{L}(\texttt{Const c})$, so `match (Const c) [] k` should raise `NoMatch`.    **Spec Satisfied!**

If `cs=c'::cs'`, then break into two further cases: either `c=c'` or `c<>c'`. In the latter case, there is no prefix of `cs` in $\mathcal{L}(\texttt{Const c})$, so again `NoMatch` should be raised. **Spec Satisfied!** However, if `c=c'`, then `[c']` $\in \mathcal{L}(\texttt{Const c})$, so whatever behavior `k` (`[c']`,`cs'`) has is what behavior `match (Const c) cs k` should have.    **Spec Satisfied!**

**Cor. 16**

For any almost total `k` and any value `c:Sigma`, (`fn cs => match (Const c) cs k`) is almost total.

# 5   Plus Correctness

**Inductive Assumption (Section 5)**

For some `R1,R2 : Sigma regexp`,

$$
\texttt{match R1 cs k} \cong
\begin{cases}
\texttt{v} & \text{where } \texttt{(p,s)} \text{ is a splitting of } \texttt{cs} \text{ such} \\
& \text{that } \texttt{p} \in \mathcal{L}(\texttt{R1}) \text{ and k accepts } \texttt{(p,s)} \\
& \text{with result v.} \\
\texttt{raise NoMatch} & \text{if there is no such } \texttt{(p,s)}
\end{cases}
$$

$$
\texttt{match R2 cs k} \cong
\begin{cases}
\texttt{v} & \text{where } \texttt{(p,s)} \text{ is a splitting of } \texttt{cs} \text{ such} \\
& \text{that } \texttt{p} \in \mathcal{L}(\texttt{R2}) \text{ and k accepts } \texttt{(p,s)} \\
& \text{with result v.} \\
\texttt{raise NoMatch} & \text{if there is no such } \texttt{(p,s)}
\end{cases}
$$

for all values `cs : Sigma list` and all almost total `k`.

In particular: if `k` is almost total, then so too are (`fn cs => match R1 cs k`) and (`fn cs => match R2 cs k`).

9

**Lemma 17**

For any almost total `k`, any `cs : Sigma list`,

$$\texttt{match (Plus(R1,R2)) cs k} \cong \begin{cases} \texttt{v} & \text{where } \texttt{(p,s)} \text{ is a split-}\\ & \text{ting of } \texttt{cs} \text{ such that } \texttt{p} \in \\ & \mathcal{L}(\texttt{Plus(R1,R2)}) \text{ and } \texttt{k} \\ & \text{ accepts } \texttt{(p,s)} \text{ with re-}\\ & \text{sult } \texttt{v}. \\ \texttt{raise NoMatch} & \text{if there is no such } \texttt{(p,s)} \end{cases}$$

**Cor. 18**

For any almost total `k`, `(fn cs => match (Plus(R1,R2)) cs k)` is almost total.

# 6   Times Correctness

**Inductive Assumption (Section 6)**

For some `R1,R2 : Sigma regexp`,

$$\texttt{match R1 cs k} \cong \begin{cases} \texttt{v} & \text{where } \texttt{(p,s)} \text{ is a splitting of } \texttt{cs} \text{ such} \\ & \text{that } \texttt{p} \in \mathcal{L}(\texttt{R1}) \text{ and } \texttt{k} \text{ accepts } \texttt{(p,s)} \\ & \text{with result } \texttt{v}. \\ \texttt{raise NoMatch} & \text{if there is no such } \texttt{(p,s)} \end{cases}$$

$$\texttt{match R2 cs k} \cong \begin{cases} \texttt{v} & \text{where } \texttt{(p,s)} \text{ is a splitting of } \texttt{cs} \text{ such} \\ & \text{that } \texttt{p} \in \mathcal{L}(\texttt{R2}) \text{ and } \texttt{k} \text{ accepts } \texttt{(p,s)} \\ & \text{with result } \texttt{v}. \\ \texttt{raise NoMatch} & \text{if there is no such } \texttt{(p,s)} \end{cases}$$

for all values `cs : Sigma list` and all almost total `k`.

In particular: if `k` is almost total, then so too are `(fn cs => match R1 cs k)` and `(fn cs => match R2 cs k)`.

> **Lemma 19**
>
> For any almost total `k`, any `cs : Sigma list`,
>
> $$\texttt{match (Times(R1,R2)) cs k} \cong \begin{cases} \texttt{v} & \text{where } \texttt{(p,s)} \text{ is a split-} \\ & \text{ting of } \texttt{cs} \text{ such that } \texttt{p} \in \\ & \mathcal{L}(\texttt{Times(R1,R2)}) \text{ and} \\ & \texttt{k} \text{ accepts } \texttt{(p,s)} \text{ with re-} \\ & \text{sult } \texttt{v}. \\ \texttt{raise NoMatch} & \text{if there is no such } \texttt{(p,s)} \end{cases}$$

> **Cor. 20**
>
> For any almost total `k`, (`fn cs => match (Times(R1,R2)) cs k`) is almost total.

# 7  Star Correctness

> **Inductive Assumption (Section 7)**
>
> For some `R : Sigma regexp` and some `cs : Sigma list`
>
> $$\texttt{match R cs g} \cong \begin{cases} \texttt{v} & \text{where } \texttt{(p,s)} \text{ is a splitting of } \texttt{cs} \text{ such} \\ & \text{that } \texttt{p} \in \mathcal{L}(\texttt{R}) \text{ and } \texttt{g} \text{ accepts } \texttt{(p,s)} \\ & \text{with result } \texttt{v}. \\ \texttt{raise NoMatch} & \text{if there is no such } \texttt{(p,s)} \end{cases}$$
>
> $$\texttt{match (Star(R)) cs' g} \cong \begin{cases} \texttt{v} & \text{where } \texttt{(p,s)} \text{ is a splitting of } \texttt{cs'} \\ & \text{such that } \texttt{p} \in \mathcal{L}(\texttt{Star(R)}) \text{ and } \texttt{g} \\ & \text{accepts } \texttt{(p,s)} \text{ with result } \texttt{v}. \\ \texttt{raise NoMatch} & \text{if there is no such } \texttt{(p,s)} \end{cases}$$
>
> for all values `cs' : Sigma list` *of length less than* `cs`, and all almost total `g`.
>
> In particular: if `g` is almost total, then so too is (`fn cs => match R cs g`). And if `g` is almost total, `match (Star(R)) cs' g` for `cs'` of length less than `cs`.

**Lemma 21**

For any almost total `k`,

$$
\texttt{match (Star(R)) cs k} \cong
\begin{cases}
\texttt{v} & \text{where } (\texttt{p},\texttt{s}) \text{ is a split-} \\
& \text{ting of } \texttt{cs} \text{ such that } \texttt{p} \in \\
& \mathcal{L}(\texttt{Star(R)}) \text{ and } \texttt{k} \text{ ac-} \\
& \text{cepts } (\texttt{p},\texttt{s}) \text{ with result } \texttt{v} \\
& . \\
\texttt{raise NoMatch} & \text{if there is no such } (\texttt{p},\texttt{s})
\end{cases}
$$

**Cor. 22**

For any almost total `k`, (`fn cs => match (Star(R)) cs k`) is almost total.

# 8    `match` and `LL` Correctness

# 9    Reduction, Representation, and Printing Specification

**Value Spec**

`represent : (''S -> string) -> ''S regexp -> string`

**REQUIRES:** `toStr` is total, `R` is non-`Zero`

**ENSURES:** `represent toStr R` converts the regular expression `R` into POSIX-like regular expression syntax, using `toStr` to convert `Sigma` values into strings. This uses the convention of representing `One` as (`.{0,0}`)

**Defn. 23**

A value `R : Sigma regexp` is said to be **Zero-reduced** if either:

- `R = Zero`

- `R` does not contain any `Zero`s

**Value Spec**

`reduce : ''S regexp -> ''s regexp`

**ENSURES:** (`reduce R`) $\implies$ `R'` such that $\mathcal{L}(\texttt{R}) = \mathcal{L}(\texttt{R'})$ and `R'` is `Zero`-reduced (and also `R'` may have some superfluous `One`s removed)

**Value Spec**

```
printRep : (''S -> string) -> ''S regexp -> unit
```

**REQUIRES:** `toStr` is total, `(reduce R)` $\not\Longrightarrow$ `Zero`

**ENSURES:** `printRep toStr R` $\Longrightarrow$ `()`

**EFFECT:** Prints the value of `(represent toStr (reduce R))` into the REPL.