

15-150 Fall 2013

Stephen Brookes

Regular expressions

1 Regular expressions

A regular expression is a sequence of characters that forms a search pattern, for use in "find and replace"-like operations on strings or character lists. The concept arose in the 1950s, when the mathematician Stephen Kleene formalized the description of a regular language (and when connections were made between automata theory and regular languages). Regular expressions came into common use with the Unix operating system, being used in the text processing utilities `ed` (an editor) and `grep` (standing for "global regular expression print", a filter operation). Many programming languages provide regular expression capabilities, some built-in (e.g. Perl, Ruby, AWK, and Tcl) and others via a standard library (e.g. Java, Python and C++).

Regular expressions give us a succinct way to identify text with certain structure, involving occurrences of specific sequences of characters, or sequences with certain general properties. You can use a regular expression to search for occurrences of a specific phrase, or an occurrence of one of a set of phrases, or a sequence of phrases, or a repeatedly occurring phrase. Once you've located a suitable phrase, you may be able to delete it, or otherwise manipulate the string. For example, you might want to search through your email for messages concerning 15-150. A regular expression search could help, using the obvious search pattern. Or you might want to look for an email message that looks like a letter from your bank. Most such letters begin with the word "Dear", followed by a name, and a colon. The last line may contain "Sincerely," "Yours truly," "Love," (probably not if it's from your bank manager!) or some other standard form. Hardly any other line in any other kind of text ends with a comma. A regular expression can look for these patterns, and a successful match tells you that a file containing them is probably a letter. It's harder to do that with an ordinary text search, unless you know the exact contents of the message you want to find.

Regular expression languages all contain (at least) the following constructs (but possibly written differently!):

- a regular expression \emptyset denoting the empty set of strings, which matches nothing;
- a regular expression ϵ denoting the set containing only the empty string, which only matches the empty string;
- a "literal" regular expression c , one for each character c in the alphabet,

matching only the string consisting of just that character;

- concatenation: $R_1 \cdot R_2$ denotes the set of strings of the form $\alpha\beta$ where α matches R_1 and β matches R_2 ;
- union: $R_1 + R_2$ denotes the set of strings that match either R_1 or R_2 ;
- iteration: R^* denotes the smallest set of strings that contains the empty string and is closed under concatenation with R .

The syntax used in different implementations of regular languages may differ drastically. For example, union is sometimes written as $R_1|R_2$, and concatenation is often written as juxtaposition, R_1R_2 . You can also use parentheses in order to disambiguate between expressions. For example, $(a + b)c$ is not the same expression as $a + (bc)$.

Many implementations of regular expressions also include additional constructs such as R^+ (which behaves the same way as $R \cdot R^*$, meaning “at least one consecutive string matching R ”) and $\{R\}$ (which behaves the same way as $(R+1)$, meaning “an optional occurrence of a string matching R ”). These two are obviously derivable from concatenation, iteration, and union.

Since there is no universal standard syntax for regular expressions, we will work with a simple ML datatype whose values serve the purpose of regular expressions. We only include the basic constructs.

To avoid confusion we will deal entirely with *lists of characters* rather than using *strings*. Concatenation will then correspond to the list append operation. There are built-in ML functions

```
explode : string -> char list
implode : char list -> string
```

that can be used to convert a string into the list of its characters, and back again.

2 ML implementation

(* A type `regexp` whose values represent "regular expressions" *)

```
datatype regexp = Zero | One | Char of char
                | Plus of regexp * regexp
                | Times of regexp * regexp
                | Star of regexp
```

The set of *values of type* `regexp` is inductively generated by the constructors of this datatype:

- `Zero` and `One` are values of type `regexp`
- For every value `c` of type `char`, `Char c` is a value of type `regexp`
- If R_1 and R_2 are values of type `regexp`, so are `Plus(R_1, R_2)` and `Times(R_1, R_2)`
- If R is a value of type `regexp`, so is `Star R`.

Every value of type `regexp` is either `Zero`, `One`, `Char c` for some value $c:\text{char}$, or one of the forms `Plus(R_1, R_2)`, `Times(R_1, R_2)` or `Star(R')`, where R_1, R_2, R' are values of type `regexp`.

3 Regular languages

A regular expression R denotes a “language” $\mathcal{L}(R)$, consisting of the set of character lists that “match” the regular expression. It may be convenient to refer to a character list as a *word*; using this parlance, a language is a set of words.

Here is an inductive mathematical definition of $\mathcal{L}(R)$, for all values R of type `regexp`.

$$\begin{aligned}\mathcal{L}(\text{Zero}) &= \emptyset \\ \mathcal{L}(\text{One}) &= \{[]\} \\ \mathcal{L}(\text{Char } c) &= \{[c]\} \\ \mathcal{L}(\text{Plus}(R_1, R_2)) &= \mathcal{L}(R_1) \cup \mathcal{L}(R_2) \\ \mathcal{L}(\text{Times}(R_1, R_2)) &= \{L_1 @ L_2 \mid L_1 \in \mathcal{L}(R_1), L_2 \in \mathcal{L}(R_2)\} \\ \mathcal{L}(\text{Star } R) &= (\mathcal{L}(R))^* = \{L_1 @ \dots @ L_n \mid n \geq 0, \text{ each } L_i \in \mathcal{L}(R)\}\end{aligned}$$

The above definition is *structural*: each clause defines the language of regular expressions built using a constructor, in terms of the languages of proper sub-expressions. For instance, the clause for $\mathcal{L}(\text{Star } R)$ says that the language of **Star** R is the set of all finite concatenations of words belonging to the language of R . We've also used the so-called *Kleene star* operation $(\cdot)^*$ in this clause, interpreted as iterated concatenation of languages.

The language of **Star** R can also be characterized recursively, since it follows from the above definition that

$$\mathcal{L}(\text{Star } R) = \{[]\} \cup \{L_1 @ L_2 \mid L_1 \in \mathcal{L}(R), L_2 \in \mathcal{L}(\text{Star } R)\}.$$

Of course we are using some obvious facts about $@$ here without proof: in particular, $@$ is associative, and for all lists L , $[] @ L = L @ [] = L$.

In fact, $\mathcal{L}(\text{Star } R)$ as defined above is the smallest set S satisfying the equation

$$S = \{[]\} \cup \{L_1 @ L_2 \mid L_1 \in \mathcal{L}(R), L_2 \in S\}.$$

To show this, let S_n be $\{L_1 @ \dots @ L_n \mid \text{each } L_i \in \mathcal{L}(R)\}$, for each $n \geq 0$. So $S_0 = \{[]\}$, $S_1 = \mathcal{L}(R)$, and so on. Let $S = \bigcup_{n=0}^{\infty} S_n$. We can easily see that $S = \{[]\} \cup \{L_1 @ L_2 \mid L_1 \in \mathcal{L}(R), L_2 \in S\}$. So we could have used this recursive equation to *define* $\mathcal{L}(\text{Star } R)$ as the least set that satisfies the equation for S . This corresponds to the way we originally described the iteration construct, adjusted to deal with character lists instead of strings: R^* denotes the smallest set containing the empty character list and closed under concatenation with (the language of) R .

The main things to remember about the language of **Star** R are:

- $[]$ is in $\mathcal{L}(\text{Star } R)$;
- Every non-empty list L belonging to $\mathcal{L}(\text{Star } R)$ is expressible as $L_1 @ L_2$ where L_1 is a list belonging to $\mathcal{L}(R)$ and L_2 belongs to $\mathcal{L}(\text{Star } R)$. In fact, we can always choose the list L_1 here to be non-empty.

Examples

Let `c` be a character, i.e. a value of type `char`. ML syntax for characters uses `#"a"` for the character “a”, and so on. And `"a"` has type `string`, not `char`.

$$\begin{aligned}\mathcal{L}(\text{Star}(\text{Char } c)) &= \{ [], [c], [c, c], \dots \} \\ \mathcal{L}(\text{Star One}) &= \{ [] \} = \mathcal{L}(\text{One}) \\ \mathcal{L}(\text{Star Zero}) &= \{ [] \} \\ \mathcal{L}(\text{Plus}(\text{Char } c, \text{One})) &= \{ [c], [] \} \\ \mathcal{L}(\text{Plus}(\text{Char } c, \text{Zero})) &= \mathcal{L}(\text{Char } c) \\ \mathcal{L}(\text{Times}(\text{Char } \# \text{"a"}, \text{Char } \# \text{"b"})) &= \{ [\# \text{"a"}, \# \text{"b"}] \} \\ \mathcal{L}(\text{Times}(\text{Char } \# \text{"a"}, \text{Times}(\text{Char } \# \text{"b"}, \text{Char } \# \text{"c"}))) &= \{ [\# \text{"a"}, \# \text{"b"}, \# \text{"c"}] \}\end{aligned}$$

Every non-empty list in $\mathcal{L}(\text{Star}(\text{Char } c))$ has the form `[c]@L2` where `L2` is in $\mathcal{L}(\text{Star}(\text{Char } c))$.

The language of `Star(Times(Char # "a", Char # "b"))` consists of all alternating lists of (the same number of) a’s and b’s.

The language of `Star(Plus(Char # "a", Char # "b"))` consists of all lists of a’s and/or b’s.

4 Implementing regular languages

We want to define an ML function

```
accepts : regexp -> char list -> bool
```

such that `accepts R L` evaluates to `true` if `L` $\in \mathcal{L}(R)$ and evaluates to `false` otherwise.

It’s not going to be easy to define `accepts` recursively! After all, the language of a regular expression may be an infinite set, and there may be many ways to break up a list into chunks, so it will not be easy to define `accepts (Star R)` in terms of `accepts R`.

Instead we’ll introduce a helper function with a more general specification, for which we can find a recursive implementation. The helper function solves a problem that we call *regular expression matching*. After we introduce the idea behind this new function we’ll show that it really would be a *helpful* function.

Regular expression matching

We'll design a helper function

```
match : regexp-> char list -> (char list -> bool) -> bool
```

such that `match R L p` evaluates to `true` if `L` has a prefix belonging to $\mathcal{L}(R)$, with a suffix that satisfies the predicate `p`; and `match R L p` evaluates to `false` otherwise. It makes sense to require that the predicate `p` is *total*, i.e. that for all character lists `A`, `p(A)` evaluates to either `true` or `false`.

In the ENSURES condition below we write `lang(R)` instead of $\mathcal{L}(R)$, meaning the language denoted by `R`.

```
(* match : regexp-> char list -> (char list -> bool) -> bool *)
(* REQUIRES p is total *)
(* ENSURES (a) match R L p = true if there are lists L1 and L2
             such that L = L1@L2, p(L2)=true, L1 is in lang(R);
           (b) match R L p = false, otherwise
*)
```

Note that a `match` function that meets this specification would definitely be a helper function for defining the `accepts` function, since `match R L null` would evaluate to `true` if `L` is in the language of `R` and evaluate to `false` otherwise.¹ So we would be able to define `accepts` this way, using `match`.

Note also that a `match` function that meets this specification also has the property that whenever `p` is a total function of type `char list -> bool` and `R` is a value of type `regexp`, the function `fn L' => match R L' p` is also a total function of type `char list -> bool`.

Also note that giving the specification is just the first step towards a satisfactory code design. We will still need to *prove* that the code is correct.

¹Remember, `null L2 = true` if and only if `L2` is the empty list, and `L@[]=L`.

The match and accepts functions

The following ML function definition is designed with this specification in mind, paying attention to the definition of $\mathcal{L}(R)$ from above.

```
fun match Zero L p = false
|   match One L p = p L
|   match (Char c) L p = (case L of
                           [ ] => false
                           | x::L' => (x = c) andalso p L')
|   match (Plus(R1, R2)) L p =
        match R1 L p orelse match R2 L p
|   match (Times(R1, R2)) L p =
        match R1 L (fn L' => match R2 L' p)
|   match (Star R') L p =
        p L orelse
        match R' L (fn L' => match (Star R') L' p)

(* accepts : regexp -> char list -> bool *)
(* ENSURES  accepts R L = true iff L is in lang(R) *)

fun accepts R L = match R L null
```

Correctness?

Try some examples to see if the `match` function produces results consistent with your knowledge of regular expressions. You should see that, for all values `c:char`,

```
match (Star (Char c)) [ ] null =>* true
match (Star (Char c)) [c] null =>* true
match (Star (Char c)) [c,c] null => true
match (Star (Char c)) [c,c] (not o null) => true
```

Try some other examples. Make sure you see how the results are consistent with the specification. For instance, in the 4th example above, `[c,c]` has a split `[c]@[c]` with the prefix `[c]` in the language of `Star(Char c)` and the suffix satisfying `(not o null) [c] = true`.

So the preliminary results of testing look favorable. However, we've warned you that testing is not an adequate substitute for a complete proof.

Let's try to prove that this code is correct.

Let $P(R)$ be the statement:

For all values L and all total functions p ,

- (a) `match R L p = true` if there are $L1, L2$ such that $L = L1 @ L2$,
`p L2 = true`, and $L1$ is in $\mathcal{L}(R)$;
- (b) `match R L p = false` otherwise.

We want to prove:

Claim

For all values $R : \text{regexp}$, $P(R)$ holds.

Proof development

Since the `match` function is defined using structural induction on regular expressions, let's try for a proof by structural induction. Here is what we need to do:

- Base cases:
 - Prove that $P(\text{Zero})$ and $P(\text{One})$ hold.
 - Prove that $P(\text{Char } c)$ holds, for all values $c : \text{char}$.
- Inductive steps:
 - Show that $P(R_1)$ and $P(R_2)$ imply $P(\text{Plus}(R_1, R_2))$.
 - Show that $P(R_1)$ and $P(R_2)$ imply $P(\text{Times}(R_1, R_2))$.
 - Show that $P(R')$ implies $P(\text{Star } R')$.

Here are the details.

- (i) For `Zero`, we prove $P(\text{Zero})$.

Let L be a character list and p be total. $\mathcal{L}(\text{Zero}) = \emptyset$.

There is no way to split L to get a suitable prefix-suffix combination. By definition, `match Zero L p = false`. So $P(\text{Zero})$ holds. (In this case, property (a) of $P(\text{Zero})$ holds trivially, because there are no relevant splits of L , and (b) holds always.)

(ii) For `One`, we prove $P(\text{One})$.

Let L be a character list and p be total. $\mathcal{L}(\text{One}) = \{[]\}$.

(a) If L has a split $L1@L2$ with $L1 \in \mathcal{L}(\text{One})$ and $p\ L2 = \text{true}$, then $L1$ must be $[]$ and $L2$ must be L . By the function definition, $\text{match One } L\ p = p(L)$. So we get $\text{match One } L\ p = \text{true}$, as required.

(b) Otherwise, if L has no such split, we must have $p\ L = \text{false}$.

This is because $[]$ is in $\mathcal{L}(\text{One})$, $L = []@L$, and p is total.

So $\text{match One } L\ p = p(L) = \text{false}$, as required.

(iii) For `Char`, let c be a value of type `char`; we prove $P(\text{Char } c)$.

Let L be a character list and p be total. $\mathcal{L}(\text{Char } c) = \{[c]\}$.

(a) If L has a split $L1@L2$ with $L1 \in \mathcal{L}(\text{Char } c)$ and $p\ L2 = \text{true}$, then $L1$ must be $[c]$ and $L = c:L2$. By definition of `match` we get

$$\text{match } (\text{Char } c)\ L\ p = \text{match } (\text{Char } c)\ (c:L2)\ p = p\ L2$$

so $\text{match } (\text{Char } c)\ L\ p = \text{true}$, as required.

(b) Otherwise, if L has no such split, there are three possibilities: either L is empty, or L does not begin with c , or L begins with c but the rest of L does not satisfy p . We consider these three possibilities:

– If L is $[]$, we have

$$\text{match}(\text{Char } c)\ L\ p = \text{match } (\text{Char } c)\ []\ p = \text{false}$$

as required.

– If $L = x:L'$ for some $x:\text{char}$ and some character list L' , but x is not the same character as c , we have

$$\begin{aligned} \text{match } (\text{Char } c)\ L\ p &= \text{match } (\text{Char } c)\ (x:L')\ p \\ &= (x = c) \text{ andalso } p\ L' \\ &= \text{false} \end{aligned}$$

as required.

– If $L = c:L'$ for some character list L' with $p(L') = \text{false}$, we have

$$\text{match } (\text{Char } c)\ L\ p = \text{match } (\text{Char } c)\ (c:L')\ p = (c = c) \text{ andalso } p\ L' =$$

as required.

In all three of these possibilities, $\text{match}(\text{Char } c) \text{ L } p = \text{false}$.

So $P(\text{Char } c)$ holds.

- (iv) For $\text{Plus}(R_1, R_2)$, assume the induction hypotheses $P(R_1)$ and $P(R_2)$.
We must show $P(\text{Plus}(R_1, R_2))$. Recall that $\mathcal{L}(\text{Plus}(R_1, R_2)) = \mathcal{L}(R_1) \cup \mathcal{L}(R_2)$.
By definition,

$$\text{match}(\text{Plus}(R_1, R_2)) \text{ L } p = (\text{match } R_1 \text{ L } p) \text{ or else } (\text{match } R_2 \text{ L } p).$$

Let p be total.

- (a) Suppose L has a split $L1@L2$ with $L1$ in $\mathcal{L}(\text{Plus}(R_1, R_2))$ and $p(L2)=\text{true}$.

We must show that $\text{match}(\text{Plus}(R_1, R_2)) = \text{true}$.

Clearly, either $L1 \in \mathcal{L}(R_1)$, or $L1 \in \mathcal{L}(R_2)$.

- (1) If $L1 \in \mathcal{L}(R_1)$, then L has a split $L1@L2$ with $L1$ in $\mathcal{L}(R_1)$ and $p(L2)=\text{true}$. By induction hypothesis $P(R_1)$ (part (a)), we have $\text{match } R_1 \text{ L } p = \text{true}$. So

$$\begin{aligned} \text{match}(\text{Plus}(R_1, R_2)) \text{ L } p &= (\text{match } R_1 \text{ L } p) \text{ or else } (\text{match } R_2 \text{ L } p) \\ &= \text{true or else } (\text{match } R_2 \text{ L } p) \\ &= \text{true} \end{aligned}$$

as required.

- (2) Otherwise, if L has no split as in (1), but it has a split $L1@L2$ with $L1 \in \mathcal{L}(R_2)$ and $p(L2)=\text{true}$, we argue as follows.

By induction hypothesis $P(R_1)$ (part (b)), $\text{match } R_1 \text{ L } p = \text{false}$.

By induction hypothesis $P(R_2)$ (part (a)), $\text{match } R_2 \text{ L } p = \text{true}$.

So

$$\begin{aligned} \text{match}(\text{Plus}(R_1, R_2)) \text{ L } p &= (\text{match } R_1 \text{ L } p) \text{ or else } (\text{match } R_2 \text{ L } p) \\ &= \text{false or else true} \\ &= \text{true} \end{aligned}$$

as required.

So in both cases (1) and (2) we get $\text{match}(\text{Plus}(R_1, R_2)) \text{ L } p = \text{true}$.

We have shown that if L has a split $L1@L2$ with $L1$ in $\mathcal{L}(\text{Plus}(R_1, R_2))$ and $p(L2)=\text{true}$, then $\text{match}(\text{Plus}(R_1, R_2)) \text{ L } p = \text{true}$.

So $P(\text{Plus}(R_1, R_2))$ (part (a)) holds.

(b) Now suppose that L has no split $L1@L2$ with $L1$ in $\mathcal{L}(\text{Plus}(R_1, R_2))$ and $p(L2)=\text{true}$. Then since $\mathcal{L}(\text{Plus}(R_1, R_2)) = \mathcal{L}(R_1) \cup \mathcal{L}(R_2)$,

- (3) L has no split $L1@L2$ with $L1$ in $\mathcal{L}(R_1)$ and $p(L2)=\text{true}$; and
- (4) L has no split $L1@L2$ with $L1$ in $\mathcal{L}(R_2)$ and $p(L2)=\text{true}$.

By the induction hypotheses for R_1 and R_2 (part (b) in both cases), we have $\text{match } R_1 \ L \ p = \text{false}$ and $\text{match } R_2 \ L \ p = \text{false}$. Hence

$$\text{match } (\text{Plus}(R_1, R_2)) \ L \ p = (\text{match } R_1 \ L \ p) \text{ or else } (\text{match } R_2 \ L \ p) = \text{false}.$$

So $P(\text{Plus}(R_1, R_2))$ (part (b)) holds, as required.

- (v) For $\text{Times}(R_1, R_2)$, assume the induction hypotheses $P(R_1)$ and $P(R_2)$. We must show $P(\text{Times}(R_1, R_2))$.

Recall that $\mathcal{L}(\text{Times}(R_1, R_2)) = \{L_1 @ L_2 \mid L_1 \in \mathcal{L}(R_1), L_2 \in \mathcal{L}(R_2)\}$.

By definition,

$$\text{match } (\text{Times}(R_1, R_2)) \ L \ p = \text{match } R_1 \ L \ (\text{fn } L' \Rightarrow \text{match } R_2 \ L' \ p).$$

Let p be total. It follows from induction hypothesis $P(R_2)$ that the function value

$$(\text{fn } L' \Rightarrow \text{match } R_2 \ L' \ p)$$

is also total.

- (a) Suppose L has a split with a prefix in $\mathcal{L}(\text{Times}(R_1, R_2))$ and suffix satisfying p . Then there are lists $L1, L2, L3$ such that $L=(L1@L2)@L3$, $L1$ in $\mathcal{L}(R_1)$, $L2$ in $\mathcal{L}(R_2)$, and $p(L3)=\text{true}$.

Since $@$ is associative, we also have $L=L1@(L2@L3)$. Let $L23$ be the value of $L2@L3$. By the properties given above, $L2 \in \mathcal{L}(R_2)$ and $pL3 = \text{true}$. So by induction hypothesis $P(R_2)$ (part (a)),

$$\text{match } R_2 \ L23 \ p = \text{true}.$$

Thus, again by the properties given above, L has a split $L1@L23$ with $L1 \in \mathcal{L}(R_1)$ and $(\text{fn } L' \Rightarrow \text{match } R_2 \ L' \ p)(L23) = \text{true}$. Since this function is also total (noted above), it follows by induction hypothesis $P(R_1)$ (part (a)) that

$$\text{match } R_1 \ L \ (\text{fn } L' \Rightarrow \text{match } R_2 \ L' \ p) = \text{true}.$$

So `match (Times(R1, R2)) L p = true`, as required.

(b) Now suppose `L` does not have a split `L1@L2` with `L1 ∈ ℒ(Times(R1, R2))` and `p(L2) = true`. We need to show that

$$\text{match } R_1 \text{ L (fn L' => match } R_2 \text{ L' p)} = \text{false}.$$

By the induction hypothesis for `R2`, the function `(fn L' => match R2 L' p)` is total. So by the induction hypothesis for `R1`, the expression

$$\text{match } R_1 \text{ L (fn L' => match } R_2 \text{ L' p)}$$

evaluates either to `true` or to `false`. To show that it evaluates to `false`, let's suppose it evaluates to `true` and derive a contradiction.

Assume this expression evaluates to `true`. By the induction hypothesis for `R1`, `L` has a split `L = A@B` such that `A ∈ ℒ(R1)` and

$$(\text{fn L' => match } R_2 \text{ L' p})(B) = \text{true}.$$

But then we also have `match R2 B p = true`.

By induction hypothesis for `R2`, `B` must have a split `B = B1@B2` with `B1 ∈ ℒ(R2)` and `p(B2) = true`. So we have `L = A@(B1@B2) = (A@B1)@B2`, and `A@B1 ∈ ℒ(Times(R1, R2))`, and `p(B2) = true`. This contradicts the original assumption that `L` had no suitable split.

Assuming that the above expression evaluates to `true` leads to a contradiction, and the expression has a value of type `bool`, so it follows that the expression evaluates to `false`. So

$$\text{match } R_1 \text{ L (fn L' => match } R_2 \text{ L' p)} = \text{false}$$

as required.

(vi) For `Star R'`, assume as induction hypothesis that `P(R')` holds. We must prove `P(Star R')`. Recall that

$$\mathcal{L}(\text{Star } R') = \{[]\} \cup \{L_1 @ L_2 \mid L_1 \in \mathcal{L}(R), L_2 \in \mathcal{L}(\text{Star } R')\}$$

and that we defined

$$\begin{aligned} \text{match (Star } R') \text{ L p} = \\ \text{p L orelse match } R' \text{ L (fn L' => match (Star } R) \text{ L' p)} \end{aligned}$$

It is obvious that this definition of `match (Star R')` is recursive, and based on the way we defined `match` for `Times`. We will need to find an inductive way to reason about this piece of code. Probably the only hope is to try using induction on the length of the character list, but then how could we show that the length decrease in all recursive calls? So now we realize there's a problem! If we try to adapt the proof argument from the `Times` case we'll need to use the structural induction hypothesis for $P(R')$ at some point, but that's only justifiable if we can prove that the function value

$$(\text{fn } L' \Rightarrow \text{match}(\text{Star } R') L' p)$$

is total. Even though we'll be assuming that `p` is total, there's no way to deduce that this function value is total. In fact it won't always be total!

After further thought you may be able to find an example of an R', L and `p` in which, even though `p` is total, `match (Star R') L p` fails to terminate. This happens if `p(L)=false` and $\mathcal{L}(R')$ contains the empty list. For example, `match (Star One) [c] null` fails to terminate.

At this stage we must abandon the correctness proof and revise the function definition! At least we have some insights into why the `Star` clause is broken.

Lesson

Our function definition seemed to be working out well, and we completed almost the entire correctness proof before discovering that the case for `Star` does not work as intended. But having failed to complete the proof we have insight into why the code doesn't work and we can see how to fix it.

The match and accepts functions, revisited

There's an easy way to fix the function: insist that `match (Star R') L p` only ever makes a recursive call to `match (Star R') L' p` when `L'` is a shorter list than `L`. Here is the revised definition of `match`. Only the clause for `Star` has been changed from the original flawed definition.

```
fun match Zero L p = false
|   match One L p = p L
|   match (Char c) L p = (case L of
                             [ ] => false
                           | x::L' => (x = c) andalso p L')
|   match (Plus(R1, R2)) L p =
      match R1 L p orelse match R2 L p
|   match (Times(R1, R2)) L p =
      match R1 L (fn L2 => match R2 L2 p)
|   match (Star R') L p =
      p L orelse
      match R' L (fn L' => length L' < length L andalso match (Star R') L' p)

(* accepts : regexp -> char list -> bool *)
(* ENSURES  accepts R L = true iff L is in lang(R) *)

fun accepts R L = match R L null
```

We can prove that this function satisfies the specification above, i.e. we can prove the claim that “For all values `R:regexp`, `P(R)` holds”, by repeating the proofs of (i) through (v) from above but with a new inductive step (vi) for the revised `Star` clause.

(vi) For `Star R'`. Recall that

$$\mathcal{L}(\text{Star } R') = \{[]\} \cup \{L_1 @ L_2 \mid L_1 \in \mathcal{L}(R), L_2 \in \mathcal{L}(\text{Star } R')\}.$$

We can rewrite this language description as

$$\mathcal{L}(\text{Star } R') = \{[]\} \cup \{L_1 @ L_2 \mid L_1 \in \mathcal{L}(R), L_1 \neq [], L_2 \in \mathcal{L}(\text{Star } R')\}.$$

(Make sure you understand why this is an equivalent description.) This formula is better suited to help analyze the correctness of the revised function definition, which is:

```

match (Star R') L p =
  p L orelse
    match R' L (fn L' => length L' < length L andalso match (Star R') L' p)

```

We want to show that $P(R')$ implies $P(\text{Star } R')$.

So assume $P(R')$ as induction hypothesis. We need to prove $P(\text{Star } R')$.

Since `match (Star R')` may make recursive calls to `match (Star R')`, we need an inductive proof for $P(\text{Star}(R'))$, and remember that we can use $P(R')$ as a hypothesis. In fact, even though we are *inside* a proof by structural induction on regular expression values, we need to use a different form of induction to deal with the `Star` case. It should be fairly obvious that we can use induction on the length of L , to prove (assuming that $P(R')$ holds):

For all values $L:\text{char list}$ and all total values p ,

- (a) `match (Star R') L p = true` if there are $L1, L2$ such that $L = L1 @ L2$, $p L2 = \text{true}$, and $L1 \in \mathcal{L}(\text{Star } R')$;
- (b) `match (Star R') L p = false` otherwise.

(This is just $P(\text{Star } R')$, rewritten.)

Proof: by induction on the length of L .

– Base case: for $L = []$.

- (a) If $[]$ has a split $L1 @ L2$ with $L1 \in \mathcal{L}(\text{Star } R')$ and $p(L2) = \text{true}$, we have $L1 = []$ and $L2 = []$, and $p[] = \text{true}$. By the function definition we have

$$\text{match (Star } R') [] p = p [] \text{ orelse } (...) = \text{true}$$

as required.

- (b) If $[]$ has no split $L1 @ L2$ with $L1 \in \mathcal{L}(\text{Star } R')$ and $p(L2) = \text{true}$, then $p[]$ does not evaluate to `true`; since p is total, it follows that $p[] = \text{false}$. So we have

$$\begin{aligned}
& \text{match (Star } R') [] p \\
&= p [] \text{ orelse } (\text{match } R' [] p') \\
&= \text{match } R' [] p'
\end{aligned}$$

where p' is the function value

$(\text{fn } L' \Rightarrow \text{length } L' < \text{length } [] \text{ andalso match (Star } R') L' p).$

Note that p' is extensionally equivalent to $(\text{fn } L' \Rightarrow \text{false})$, and that p' is total. Also there is no split of $[]$ with a suffix that satisfies p' . So by hypothesis $P(R')$, $\text{match } R' [] p' = \text{false}$. So we've shown that

$$\text{match (Star } R') [] p = \text{false},$$

as required.

- Inductive step: Let L be a non-empty list, and assume that for all lists L' shorter than L ,

(a') $\text{match (Star } R') L' p = \text{true}$ if there are $L1, L2$ such that $L' = L1 @ L2$, $p L2 = \text{true}$, and $L1 \in \mathcal{L}(\text{Star } R')$;

(b') $\text{match (Star } R') L' p = \text{false}$ otherwise.

Need to show (a) and (b) for L .

- (a) Suppose L has a split $L1 @ L2$ with $L1 \in \mathcal{L}(\text{Star } R')$ and $p(L2) = \text{true}$.

If $p L = \text{true}$ we can use the split $[] @ L$ of L as in the base case, so again we get $\text{match (Star } R') L p = \text{true}$ as required. Otherwise, we have $p L = \text{false}$ and we have a split with $L1$ non-empty, $L2$ shorter than L , satisfying the above properties. Hence

$$\begin{aligned} \text{match (Star } R') L p \\ &= p L \text{ orelse } (\text{match } R' L p') \\ &= \text{match } R' L p' \end{aligned}$$

where p' is the function value

$(\text{fn } L' \Rightarrow \text{length } L' < \text{length } L \text{ andalso match (Star } R') L' p).$

By the induction hypotheses (a') and (b'), p' is a total function. But $L1$ is a non-empty list belonging to $\mathcal{L}(\text{Star } R')$, so must have a split $L1 = A @ B$ with A non-empty, $A \in \mathcal{L}(R')$ and $B \in \mathcal{L}(\text{Star } R')$.

Let $B2$ be the value of $B@L2$. Then $\text{length } B2 < \text{length } L$, so

$$p'(B@L2) = p'(B2) = \text{match } (\text{Star } R') \ B2 \ p.$$

But, from above, we know that the list $B2$ is shorter than L and has a split $B@L2$ with $B \in \mathcal{L}(\text{Star } R')$ and $p(L2) = \text{true}$. So by induction hypothesis (a'), $\text{match } (\text{Star } R') \ B2 \ p = \text{true}$. Hence $p'(B2) = \text{true}$. So L has a split $A@B2$ with $A \in \mathcal{L}(R')$ and $p'(B2) = \text{true}$. It follows from $P(R')$ that $\text{match } (\text{Star } R') \ L \ p = \text{true}$, as required.

- (b) We leave the details for this part as an exercise. Try adapting the argument given above for the case where L is $[\]$.