


Lec 14 RegEx

Reg Ex

* Chomsky ... ? language hierarchy

Ex.  What we can get:

- c	- c*
- cc	- cgc
- ccc	...

Operations: e.g. $(g+w)^*c$ ^{or}

Notation & Definition

Σ - alphabet of characters

Σ^* - all finite string over Σ $aabba \in \{a,b\}^*$

ϵ - empty string

Language - subset of Σ^*

Regular expression:

a for all $a \in \Sigma$
0
1 } special symbols

$r_1 + r_2$ or
 $r_1 r_2$ concat
 r^* repeat } with r, r_1, r_2 being RegEx

Regular Language \leftrightarrow Finite automata, tokenisation

SML is a context-free language
Turing machine is regular plus
two stacks

Regular language denoted $L(r)$

$$L(a) = \{a\} \leftarrow \text{singleton}$$

$$L(\emptyset) = \{\} \leftarrow \text{empty language}$$

$$L(\epsilon) = \{\epsilon\} \leftarrow \text{language with empty string only} \quad L(\epsilon) = \{\epsilon\}$$

$$L(r_1 + r_2) = \{s \mid s \in L(r_1) \vee s \in L(r_2)\}$$

$$L(r_1 r_2) = \{s \mid s \in L(r_1) \wedge s \in L(r_2)\}$$

$$L(r^*) = \{s \mid s = s_1 s_2 \dots s_n, n \geq 0, \text{ each } s_i \in L(r)\} \quad \text{note } \epsilon \in L(r^*)$$

* L is regular if $L = L(r)$ f.s. regex r .

Ex. suppose $\Sigma = \{a, b\}$

$$- L(aa) = \{aa\}$$

$$- L((a+b)^*) = \Sigma^*$$

$$- L((a+b)^* aa (a+b)^*) = \text{all strs in } \Sigma^* \text{ that has "aa" in it}$$

$$- L((a+1)(b+ba)^*) = \text{all strs in } \Sigma^* \text{ without "aa" in it}$$

$$- \underline{L(r^*) = L(1+rr^*)}$$

Notice multiple can mean same thing

Acceptor

(* accept : regexp \rightarrow string \rightarrow bool

REQ true

ENS accept r $s \mapsto$ true if $s \in L(r)$
 \mapsto false otherwise

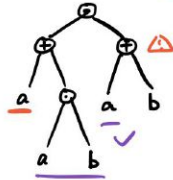
*)

Consider: $r = (a+ab)(a+b) \rightarrow L(r) = \{aa, ab, aba, abb\}$

$\frac{ab}{a}$ $\frac{a}{ba}$

notice we may not know where to split. Need backtracking \rightarrow Continuation!!

There's actually backtracking going on. See aside trace



* match: regexp \rightarrow char list \rightarrow (char list \rightarrow bool) \rightarrow bool

REQ k total

ENS match r cs k \hookrightarrow true if $cs \cong p@s$ s.t. $p \in L(r) \ \& \ k(cs) \cong$ true
 \hookrightarrow false otherwise

*)

Code

turn into char list check if empty

fun accept r s = match r (String.explode s) List.null

datatype regexp =
 | Zero
 | One
 | Plus of regexp * regexp
 | Times of regexp * regexp
 | Star of regexp

fun match (Char a) cs k =
 (case cs of
 [] \Rightarrow false | \leftarrow can't split out p
 c::cs' \Rightarrow (a=c) andalso (k cs')

| match Zero -- = false \leftarrow nothing in that language

$| \text{match One } cs \ k = k \ cs$
 $| \text{match (Plus } (r_1, r_2)) \ cs \ k =$
 $\quad (\text{match } r_1 \ cs \ k) \ \text{or else } (\text{match } r_2 \ cs \ k)$
 $| \text{match (Times } (r_1, r_2)) \ cs \ k =$
 $\quad \text{match } r_1 \ cs \ (\text{fn } cs' \Rightarrow \text{match } r_2 \ cs' \ k)$
 $| \text{match (Star } (r)) \ cs \ k = \leftarrow \text{ has bug} \quad \text{Recall: } L(r^*) = L(1+rr^*)$
 $\quad (k \ cs) \ \text{or else } (\text{match } r \ cs \ (\text{fn } cs' \Rightarrow \text{match (Star } r) \ cs' \ k))$

Proof-directed Debugging

We can't prove $(\text{fn } cs' \Rightarrow \text{match (Star } r) \ cs' \ k)$ is total
 Try: $\text{match (Star One) } [#"a"] \ \text{List.null}$, that loops!

Solution ?

1. Require input regex in standard form

Standard form \Leftrightarrow for any subexpression $\text{Star}(r')$, $\epsilon \notin r'$ \leftarrow looping to look for empty is bad: C
 We can actually write preprocessor to get any r to this form

2. Runtime check that cs' is shorter than cs , viz. cs' is proper prefix

$| \text{match (Star } (r)) \ cs \ k =$
 $\quad (k \ cs) \ \text{or else } (\text{match } r \ cs \ (\text{fn } cs' \Rightarrow$
 $\quad \quad \text{properSuffix } (cs', cs)$
 $\quad \quad \text{and also}$
 $\quad \quad \text{match (Star } r) \ cs' \ k)$
 $\quad)$

Aside: example trace

Let $r = (a+ab)(a+ab)$ (actually array)

accept r "aba" 

\Rightarrow match r "aba" List.null

\Rightarrow match $(a+ab)$ "aba" (fn cs' \Rightarrow match $(a+ab)$ cs' List.null)

\Rightarrow (match a "aba" k_1) or else (match ab "aba" k_1)

\Rightarrow k_1 , "ba" \sim

\Rightarrow match $(a+ab)$ "ba" List.null \sim

$\dots \Rightarrow$ false \sim

\cong (match ab "aba" k)

$\dots \Rightarrow$ true

Proof of Correctness

Assume termination

WTS $\varphi(r) \equiv \text{match } r \text{ cs } k \cong \text{true} \iff \exists p, s, \text{cs} \cong p@s, p \in L(r), k(\text{cs}) \cong \text{true}$

\Rightarrow completeness

\Leftarrow soundness

Go structural induction on r

BC1: Zero

BC2: One [omit]

BC3: Char (a)

IC1: $r = \text{Plus}(r_1, r_2)$

IH: $\varphi(r_1), \varphi(r_2)$

Soundness Suppose $\text{match}(\text{Plus}(r_1, r_2)) \text{ cs } k \cong \text{true}$

LHS $\cong (\text{match } r_1 \text{ cs } k) \text{ or else } (\text{match } r_2 \text{ cs } k)$

By IH, $\exists p, s, p@s \in L(r_1), k \text{ s true}$

Then $p \in L(\text{Plus}(r_1, r_2))$ as req

Completeness Suppose $\exists p, s, \text{cs} \cong p@s, p \in L(\text{Plus}(r_1, r_2)), k(\text{cs}) \cong \text{true}$

} [omit]
v