

Introduction to Regular Expression Derivatives

A functional, algebraic approach to generate DFAs

Schrodinger Yifan ZHU

Computer Science Department
University of Rochester

Agenda

- ▶ Derivative-based DFA generation for individual regular expressions
 - ▶ Regular expression definitions (set-based)
 - ▶ Derivatives of REs
 - ▶ Weak-equivalence between REs
 - ▶ Congruence classes of the alphabet
- ▶ Derivative-based DFA generation for lexers (groups of REs)



What are regular expressions?

Regular expressions on an ambient set Σ (a.k.a. alphabet) is inductively defined with the following rules:

$$\mathcal{L} := \mathcal{S} \mid \varepsilon \mid \alpha \sim \beta \mid \alpha + \beta \mid \alpha \& \beta \mid \alpha^* \mid \neg \alpha \mid \emptyset \quad (1)$$

- ▶ We may abuse single characters ('a') to represent singleton sets.
- ▶ We use some pictures from Scott Owens' Regular-expression derivatives re-examined, where \sim may be omitted or replaced by \cdot instead.



What are regular expressions?

Regular expressions on an ambient set Σ (a.k.a. alphabet) is inductively defined with the following rules:

- ▶ $\mathcal{S} \subseteq \Sigma$: single character languages where the character comes from a subset of the alphabet.
- ▶ ϵ : empty languages.
- ▶ $\alpha \sim \beta$: languages where sub-language α immediately followed by sub-language β .
- ▶ $\alpha + \beta$: union of languages represented by α or β .
- ▶ $\alpha \& \beta$: intersection of languages represented by α and β .
- ▶ α^* : Kleene closure of α .
- ▶ \emptyset : no language.
- ▶ $\neg\alpha$: languages that α rejects.



What are regular expression derivatives?

A **derivative** of a regular expression \mathcal{L} over a character is another regular expression \mathcal{L}' that accepts languages that are accepted in \mathcal{L} after character is consumed.

$$\partial_a(a \sim b \sim c) = b \sim c$$

$$\partial_b(a \sim b \sim c) = \emptyset$$

$$\partial_a((a + b) \sim b \sim c) = b \sim c$$



What are regular expression derivatives?

$$\partial_{\gamma}(\mathcal{S}) = \begin{cases} \varepsilon, & \gamma \in \mathcal{S} \\ \emptyset, & \gamma \notin \mathcal{S} \end{cases}$$

$$\partial_{\gamma}(\varepsilon) = \emptyset$$

$$\partial_{\gamma}(\alpha \sim \beta) = \begin{cases} (\partial_{\gamma}(\alpha) \sim \beta) + \partial_{\gamma}(\beta), & \text{nullable}(\alpha) \\ \partial_{\gamma}(\alpha) \sim \beta, & \neg \text{nullable}(\alpha) \end{cases} \quad (2)$$

$$\partial_{\gamma}(\alpha + \beta) = \partial_{\gamma}(\alpha) + \partial_{\gamma}(\beta)$$

$$\partial_{\gamma}(\alpha \& \beta) = \partial_{\gamma}(\alpha) \& \partial_{\gamma}(\beta)$$

$$\partial_{\gamma}(\alpha^*) = \partial_{\gamma}(\alpha) \sim \alpha^*$$

$$\partial_{\gamma}(\neg \alpha) = \neg \partial_{\gamma}(\alpha)$$



What are regular expression derivatives?

nullable (S) = **False**

nullable (\emptyset) = **False**

nullable (ϵ) = **True**

nullable ($\alpha \sim \beta$) = **nullable** (α) \wedge **nullable** (β) (3)

nullable ($\alpha + \beta$) = **nullable** (α) \vee **nullable** (β)

nullable ($\alpha\&\beta$) = **nullable** (α) \wedge **nullable** (β)

nullable (α^*) = **True**

nullable ($\neg\alpha$) = \neg **nullable** (α)



The big picture

- ▶ Each regular expression \mathcal{L} represents a state in DFA.
- ▶ At each state \mathcal{L} , enumerate the alphabet Σ and create a edge from \mathcal{L} to $\partial_\gamma(\mathcal{L})$ for character γ .
- ▶ Successful matches are represented by **nullable** expressions.
- ▶ Rejecting states are indicated by \emptyset .



The big picture

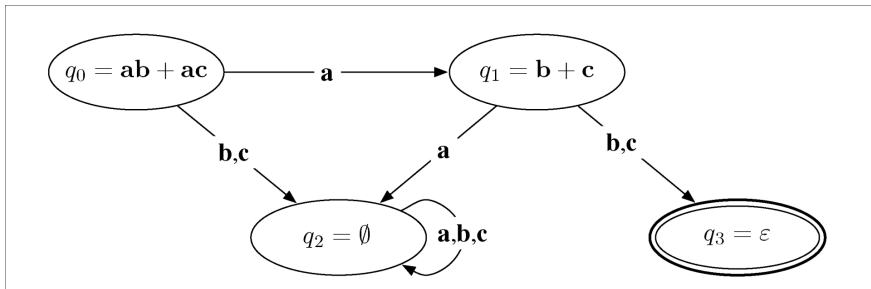


Figure: An example DFA for $(a \sim b) + (a \sim c)$

Semantically identical regular expressions?

For regular expressions with intersection and complement, checking if two expressions accept the same set of languages may acquire **non-elementary time complexity** (Aho, et al. The Design and Analysis of Computer Algorithms)



A weaker notion of equivalence

$$\begin{array}{ll} r \& r \approx r & (*) \quad r + r \approx r \\ r \& s \approx s \& r & (*) \quad r + s \approx s + r \\ (r \& s) \& t \approx r \& (s \& t) & (*) \quad (r + s) + t \approx r + (s + t) \\ \emptyset \& r \approx \emptyset & \neg \emptyset + r \approx \neg \emptyset \\ \neg \emptyset \& r \approx r & \emptyset + r \approx r \\ \\ (r \cdot s) \cdot t \approx r \cdot (s \cdot t) & (r^*)^* \approx r^* \\ \emptyset \cdot r \approx \emptyset & \varepsilon^* \approx \varepsilon \\ r \cdot \emptyset \approx \emptyset & \emptyset^* \approx \varepsilon \\ \varepsilon \cdot r \approx r & \neg(\neg r) \approx r \\ r \cdot \varepsilon \approx r & \end{array}$$

Figure: A weaker notion of equivalence between REs.



A weaker notion of equivalence

- ▶ One can add/remove rules based on practical cases. (DotNet's alternation rule is not commutative)
- ▶ If you see the rule from left to right, it actually provides a way to simplify and reorder rules.
- ▶ We use **Wnorm**(\mathcal{L}) to represent apply weakly normalization on expression \mathcal{L} based on this weaker notion of equivalence.



Congruence classes

Iterate through the whole alphabet?
Too costly for larger alphabets!



Congruence classes

Given two characters $\alpha, \beta \in \Sigma$ and a regular expression \mathcal{L} (denoted by $\alpha \simeq_{\mathcal{L}} \beta$), we say α is **congruent** to β at expression \mathcal{L} if $\partial_{\alpha}(\mathcal{L})$ and $\partial_{\beta}(\mathcal{L})$ accept the same set of languages.

$$a \simeq_{(a|b) \sim c} b \iff \partial_a((a|b) \sim c) \equiv \partial_b((a|b) \sim c) \equiv c$$

$$d \simeq_{(a|b) \sim c} e \iff \partial_d((a|b) \sim c) \equiv \partial_e((a|b) \sim c) \equiv \emptyset$$



Congruence classes

- ▶ If we know all the congruence classes for \mathcal{L} , we only need to calculate the derivatives per class.
- ▶ Still no precise solution within reasonable time complexity.
- ▶ Finer partitions are safe!



Congruence classes

- ▶ If we know all the congruence classes for \mathcal{L} , we only need to calculate the derivatives per class.
- ▶ Still no precise solution within reasonable time complexity.
- ▶ Finer partitions are safe!



Congruence classes

$$\mathbf{ApproxCC}(\epsilon) = \{\Sigma\}$$

$$\mathbf{ApproxCC}(\mathcal{S}) = \{\mathcal{S}, \Sigma - \mathcal{S}\}$$

$$\mathbf{ApproxCC}(\alpha \sim \beta) = \begin{cases} \mathbf{ApproxCC}(\alpha), & \neg \mathbf{nullable}(\alpha) \\ \mathbf{ApproxCC}(\alpha) \wedge \mathbf{ApproxCC}(\beta), & \mathbf{nullable}(\alpha) \end{cases}$$

$$\mathbf{ApproxCC}(\alpha + \beta) = \mathbf{ApproxCC}(\alpha) \wedge \mathbf{ApproxCC}(\beta)$$

$$\mathbf{ApproxCC}(\alpha \& \beta) = \mathbf{ApproxCC}(\alpha) \wedge \mathbf{ApproxCC}(\beta)$$

$$\mathbf{ApproxCC}(\alpha^*) = \mathbf{ApproxCC}(\alpha)$$

$$\mathbf{ApproxCC}(\neg \alpha) = \mathbf{ApproxCC}(\alpha)$$

$$S \wedge T = \{s \cap t \mid s \in S \vee t \in T\}$$

(4)



Congruence classes

$$\begin{aligned} & \mathbf{ApproxCC}((a + (b \sim a)) \sim c) \\ = & \mathbf{ApproxCC}(a + (b \sim a)) \wedge \mathbf{ApproxCC}(c) \\ = & \mathbf{ApproxCC}(a) \wedge \mathbf{ApproxCC}(b \sim a) \wedge \mathbf{ApproxCC}(c) \\ = & \mathbf{ApproxCC}(a) \wedge \mathbf{ApproxCC}(b) \wedge \mathbf{ApproxCC}(c) \\ = & \{a, \Sigma - a\} \wedge \{b, \Sigma - b\} \wedge \{c, \Sigma - c\} \\ = & \{\emptyset, a, b, \Sigma - \{a, b\}\} \wedge \{c, \Sigma - c\} \\ = & \{\emptyset, a, b, c, \Sigma - \{a, b, c\}\} \end{aligned}$$

Precise result is $\{\{a, c\}, b, \Sigma - \{a, b, c\}\}$.



DFA generation

Algorithm 1 DFA states exploration

```
procedure EXPLOREDFA( $x, V, T$ )  $\triangleright x$  is current state (normalized),  $V$  is the set of visited states  
  if  $x' \in V$  then  
    return  
  end if  
   $V \leftarrow V \cup \{x'\}$   
  for  $c \in \text{APPROXCC}(x')$  do  
    if  $\text{ISEMPTY}(c)$  then  
      continue  
    end if  
     $r \leftarrow \text{REPRESENTATIVE}(c)$   
     $y \leftarrow \text{WNORM}(\partial_r(x'))$   
     $T[x'][c] \leftarrow y$   $\triangleright$  update DFA transition table  
    EXPLOREDFA( $y, V, T$ )  
  end for  
end procedure
```

Figure: DFA States Exploration



Regular expression vectors for lexer generation

A **regular expression vector** $(l_1, l_2, \dots, l_n) \in \mathcal{L}^n$ is just a tuple of regular expressions.

$$\mathbf{Wnorm}((l_1, \dots, l_n)) = (\mathbf{Wnorm}(l_1), \dots, \mathbf{Wnorm}(l_n))$$

$$\mathbf{ApproxCC}((l_1, \dots, l_n)) = \bigwedge_{i=1}^n \mathbf{ApproxCC}(l_i) \quad (5)$$

$$\partial((l_1, \dots, l_n)) = (\partial(l_1), \dots, \partial(l_n))$$

(\bigwedge is the meet operator defined in **ApproxCC**.)



Regular expression vectors for lexer generation

$$\begin{aligned}\mathbf{Accept}((l_1, l_2, \dots, l_n)) &= \bigvee_{i=1}^n \mathbf{nullable}(l_i) \\ \mathbf{Reject}((l_1, l_2, \dots, l_n)) &= \bigwedge_{i=1}^n l_i \equiv \emptyset\end{aligned}\tag{6}$$

(\bigwedge, \bigvee are logical operators.)



Longest-match lexer

Algorithm 2 Longest-match Lexer

```
procedure LEXER(S) ▷ S is the input
  M ← 0 ▷ M is the longest match, initialized to 0 (no match)
  L ← 0 ▷ L is offset of the longest match
  l ← LENGTH(S)
  s ← State0
  for i ∈ RANGE(0, l) do
    c ← S[i] ▷ c is set to EOF if the lexer is at the end of the input
    if s = State0 then
      if comptime ACCEPT(State0) then
        L ← i
        M ← comptime FIRSTACCEPT(State0)
      end if
      if comptime REJECT(State0) then
        break
      end if
      if c ∈ comptime APPROXCC(State0)[0] then
        s ← comptime T[State0][APPROXCC(State0)[0]]
      end if
      go though all congruence classes
    end if
    go though all states
  end for
  return (M, L)
end procedure
```



Num of States

Lexer	ml-lex	ml-ulex	Minimal	Description
Burg	61	58	58	A tree-pattern match generator
CKit	122	115	115	ANSI C lexer
Calc	12	12	12	Simple calculator
CM	153	146	146	The SML/NJ compilation manager
Expression	19	19	19	A simple expression language
FIG	150	144	144	A foreign-interface generator
FOL	41	41	41	First-order logic
HTML	52	49	49	HTML 3.2
MDL	161	158	158	A machine-description language
ml-lex	121	116	116	The ml-lex lexer
Scheme	324	194	194	R ⁵ RS Scheme
SML	251	244	244	Standard ML lexer
SML/NJ	169	158	158	SML/NJ lexer
Pascal	60	55	55	Pascal lexer
ml-yacc	100	94	94	The ml-yacc lexer
Russo	4803	3017	2892	System-log data mining
L ₂	<i>n/a</i>	147	106	Monitoring stress-test



Preprocessing Time

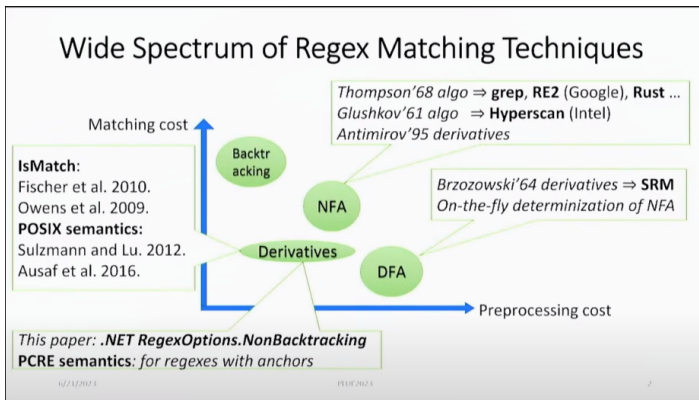


Figure: PLDI slides for Derivative Based Nonbacktracking Real-World Regex Matching with Backtracking Semantics