

Symbolic Solving of Extended Regular Expression Inequalities

Matthias Keil, Peter Thiemann
University of Freiburg,
Freiburg, Germany

December 15, 2014, IARCS Annual Conference on Foundations of
Software Technology and Theoretical Computer Science



Notizen

Extended Regular Expressions

Definition

$$r, s, t := \epsilon \mid A \mid r+s \mid r \cdot s \mid r^* \mid r \& s \mid !r$$

- Σ is a potentially infinite set of symbols
- $A, B, C \subseteq \Sigma$ range over sets of symbols
- $\llbracket r \rrbracket \subseteq \Sigma^*$ is the language of a regular expression r ,
where $\llbracket A \rrbracket = A$

Notizen

Language Inclusion

Definition

Given two regular expressions r and s ,

$$r \sqsubseteq s \Leftrightarrow \llbracket r \rrbracket \subseteq \llbracket s \rrbracket$$

- $\llbracket r \rrbracket \subseteq \llbracket s \rrbracket$ iff $\llbracket r \rrbracket \cap \overline{\llbracket s \rrbracket} = \emptyset$
- Decidable using standard techniques:
Construct DFA for $r \& !s$ and check for emptiness
- Drawback is the expensive construction of the automaton
- PSPACE-complete

Notizen

Antimirov's Algorithm



- Deciding containment for *basic regular expressions*
- Based on derivatives and expression rewriting
- Avoid the construction of an automaton
- $\partial_a(r)$ computes a regular expression for $a^{-1}[[r]]$ (Brzozowski) with $u \in [[r]]$ iff $\epsilon \in [[\partial_u(r)]]$

Lemma

For regular expressions r and s ,

$$r \sqsubseteq s \Leftrightarrow (\forall u \in \Sigma^*) \partial_u(r) \sqsubseteq \partial_u(s).$$

Notizen

Antimirov's Algorithm (cont'd)



Lemma

$$r \sqsubseteq s \Leftrightarrow (\nu(r) \Rightarrow \nu(s)) \wedge (\forall a \in \Sigma) \partial_a(r) \sqsubseteq \partial_a(s)$$

$$\frac{\text{CC-DISPROVE} \quad \nu(r) \wedge \neg \nu(s)}{r \sqsubseteq s \vdash_{\text{CC}} \text{false}}$$

$$\frac{\text{CC-UNFOLD} \quad \nu(r) \Rightarrow \nu(s)}{r \sqsubseteq s \vdash_{\text{CC}} \{\partial_a(r) \sqsubseteq \partial_a(s) \mid a \in \Sigma\}}$$

- Choice of next step's inequality is nondeterministic
- An infinite alphabet requires to compute for infinitely many $a \in \Sigma$

Notizen

First Symbols



Lemma

$$r \sqsubseteq s \Leftrightarrow (\nu(r) \Rightarrow \nu(s)) \wedge (\forall a \in \text{first}(r)) \partial_a(r) \sqsubseteq \partial_a(s)$$

- Let $\text{first}(r) := \{a \mid a\omega \in [[r]]\}$ be the set of first symbols
- Restrict symbols to first symbols of the left hand side
- CC-UNFOLD does not have to consider the entire alphabet
- For extended regular expressions, $\text{first}(r)$ may still be an infinite set of symbols

Notizen

Problems



- Antimirov's algorithm only works with basic regular expressions or requires a finite alphabet
- Extension of *partial derivatives* (Caron et al.) that computes an NFA from an extended regular expression
- Works on sets of sets of expressions
- Computing derivatives becomes more expensive

Notizen

Goal



- Algorithm for deciding $[r] \subseteq [s]$ quickly
- Handle *extended regular expressions*
- Deal effectively with very large (or infinite) alphabets (e.g. Unicode character set)

Solution

- Require finitely many atoms, even if the alphabet is infinite
- Compute derivatives with respect to literals

Notizen

Representing Sets of Symbols



A literal is a set of symbols $A \subseteq \Sigma$

Definition

A is an element of an *effective* boolean algebra $(U, \sqcup, \sqcap, \bar{\cdot}, \perp, \top)$ where $U \subseteq \wp(\Sigma)$ is closed under the boolean operations.

- For finite (small) alphabets:
 $U = \wp(\Sigma)$, $A \subseteq \Sigma$
- For infinite (or just too large) alphabets:
 $U = \{A \in \wp(\Sigma) \mid A \text{ finite} \vee \bar{A} \text{ finite}\}$
- Second-level regular expressions:
 $\Sigma \subseteq \wp(\Gamma^*)$ with $U = \{A \subseteq \wp(\Gamma^*) \mid A \text{ is regular}\}$
- Formulas drawn from a first-order theory over alphabets
For example, $[a-z]$ represented by $x \geq 'a' \wedge x \leq 'z'$

Notizen

Derivatives with respect to Literals



- Definition for $\partial_A(r)$?
- $\partial_a(r)$ computes a regular expression for $a^{-1}[[r]]$ (Brzozowski)

Desired property

$$[[\partial_A(r)]] \stackrel{?}{=} A^{-1}[[r]] = \bigcup_{a \in A} a^{-1}[[r]] = \bigcup_{a \in A} [[\partial_a(r)]]$$

Notizen

Positive Derivatives on Literals



Definition

$$\delta_A^+(B) := \begin{cases} \epsilon, & B \cap A \neq \perp \\ \emptyset, & \text{otherwise} \end{cases}$$

Problem

With $A = \{a, b\}$ and $r = (a \cdot c) \& (b \cdot c)$,

$$\begin{aligned} \delta_A^+(r) &= \delta_A^+(a \cdot c) \& \delta_A^+(b \cdot c) \\ &= c \& c \\ &\sqsubseteq \emptyset \end{aligned}$$

Notizen

Negative Derivatives on Literals



Definition

$$\delta_A^-(B) := \begin{cases} \epsilon, & \overline{B} \cap A = \perp \\ \emptyset, & \text{otherwise} \end{cases}$$

Problem

With $A = \{a, b\}$ and $r = (a \cdot c) + (b \cdot c)$,

$$\begin{aligned} \delta_A^-(r) &= \delta_A^-(a \cdot c) + \delta_A^-(b \cdot c) \\ &= \emptyset + \emptyset \\ &\sqsubseteq c \end{aligned}$$

Notizen

Positive and Negative Derivatives



- Extends Brzozowski's derivative operator to sets of symbols.
- Defined by induction and flip on the complement operator

Definition

From $\partial_a(!s) = !\partial_a(s)$, define:

$$\delta_A^+(!r) := !\delta_A^-(r) \quad | \quad \delta_A^-(!r) := !\delta_A^+(r)$$

Lemma

For any regular expression r and literal A ,

$$\llbracket \delta_A^+(r) \rrbracket \supseteq \bigcup_{a \in A} \llbracket \partial_a(r) \rrbracket \quad | \quad \llbracket \delta_A^-(r) \rrbracket \subseteq \bigcap_{a \in A} \llbracket \partial_a(r) \rrbracket$$

Notizen

Literals of an Inequality



Lemma

$$r \sqsubseteq s \Leftrightarrow (\nu(r) \Rightarrow \nu(s)) \wedge (\forall a \in \text{first}(r)) \partial_a(r) \sqsubseteq \partial_a(s)$$

- $\text{first}(r)$ may still be an infinite set of symbols
- Use *first literals* as representatives of the *first symbols*

Example

- 1 Let $r = \{a, b, c, d\} \cdot d^*$, then $\{a, b, c, d\}$ is a first literal
- 2 Let $s = \{a, b, c\} \cdot c^* + \{b, c, d\} \cdot d^*$, then $\{a, b, c\}$ and $\{b, c, d\}$ are first literals

Notizen

Literals of an Inequality (cont'd)



Problem

Let $r = \{a, b, c, d\} \cdot d^*$, $s = \{a, b, c\} \cdot c^* + \{b, c, d\} \cdot d^*$, and $A = \{a, b, c, d\}$, then

$$\begin{aligned} \delta_A^+(r) &\sqsubseteq \delta_A^+(s) & (1) \\ \delta_A^+(\{a, b, c, d\} \cdot d^*) &\sqsubseteq \delta_A^+(\{a, b, c\} \cdot c^*) + \delta_A^+(\{b, c, d\} \cdot d^*) & (2) \\ d^* &\sqsubseteq c^* + d^* & (3) \end{aligned}$$

- Positive (negative) derivatives yield an upper (lower) approximation
- To obtain the precise information, we need to restrict these literals suitably to *next literals*, e.g. $\{\{a\}, \{b, c\}, \{d\}\}$

Notizen

Next Literals



$$\begin{aligned} \text{next}(\epsilon) &= \{\emptyset\} \\ \text{next}(A) &= \{A\} \\ \text{next}(r+s) &= \text{next}(r) \times \text{next}(s) \\ \text{next}(r \cdot s) &= \begin{cases} \text{next}(r) \times \text{next}(s), & \nu(r) \\ \text{next}(r), & \neg \nu(r) \end{cases} \\ \text{next}(r^*) &= \text{next}(r) \\ \text{next}(r \&s) &= \text{next}(r) \sqcap \text{next}(s) \\ \text{next}(!r) &= \text{next}(r) \cup \{\sqcap \{ \bar{A} \mid A \in \text{next}(r) \}\} \end{aligned}$$

Definition

Let \mathcal{L}_1 and \mathcal{L}_2 be two sets of disjoint literals.

$$\mathcal{L}_1 \times \mathcal{L}_2 := \{(A_1 \sqcap A_2), (A_1 \sqcap \overline{\sqcup \mathcal{L}_2}), (\overline{\sqcup \mathcal{L}_1} \sqcap A_2) \mid A_1 \in \mathcal{L}_1, A_2 \in \mathcal{L}_2\}$$

Notizen

Next Literals (cont'd)



Example

Let $s = \{a, b, c\} \cdot c^* + \{b, c, d\} \cdot d^*$, then

$$\begin{aligned} \text{next}(s) &= \text{next}(\{a, b, c\} \cdot c^* + \{b, c, d\} \cdot d^*) \\ &= \{\{a, b, c\}\} \times \{\{c\}\} + \{\{b, c, d\}\} \times \{\{d\}\} \\ &= \{\{a\}, \{b, c\}, \{d\}\} \end{aligned}$$

Lemma

For all r ,

- $\bigcup \text{next}(r) \supseteq \text{first}(r)$
- $|\text{next}(r)|$ is finite
- $(\forall A, B \in \text{next}(r)) A \sqcap B = \emptyset$

Notizen

Coverage



Lemma

Let $\mathcal{L} = \text{next}(r)$ and $A \in \text{next}(r) \setminus \{\emptyset\}$.

- 1 $(\forall a, b \in A) \partial_a(r) = \partial_b(r) \wedge \delta_A^+(r) = \delta_{\bar{A}}(r) = \partial_a(r)$
- 2 $(\forall a \notin \bigcup \mathcal{L}) \partial_a(r) = \emptyset$

Definition

Let $A' \in \text{next}(r)$. For each $\emptyset \neq A \subseteq A'$ define $\partial_A(r) := \partial_a(r)$, where $a \in A$.

Notizen

Next Literals of an Inequality



- Next literal of $\text{next}(r \sqsubseteq s)$
- Sound to join literals of both sides $\text{next}(r) \times \text{next}(s)$
- Contains also symbols from s
- First symbols of r are sufficient to prove containment

Definition

Let \mathcal{L}_1 and \mathcal{L}_2 be two sets of disjoint literals.

$$\mathcal{L}_1 \times \mathcal{L}_2 := \{(A_1 \sqcap A_2), (A_1 \sqcap \overline{\sqcap \mathcal{L}_2}) \mid A_1 \in \mathcal{L}_1, A_2 \in \mathcal{L}_2\}$$

Left-based join corresponds to $\text{next}(r \&!s)$.

Definition

Let $r \sqsubseteq s$ be an inequality, define: $\text{next}(r \sqsubseteq s) := \text{next}(r) \times \text{next}(s)$

Notizen

Solving Inequalities



Lemma

$$r \sqsubseteq s \Leftrightarrow (\nu(r) \Rightarrow \nu(s)) \wedge (\forall a \in \text{first}(r)) \partial_a(r) \sqsubseteq \partial_a(s)$$

To determine a finite set of representatives

- select *one* symbol a from each equivalence class $A \in \text{next}(r)$
- calculate with $\delta_A^+(r)$ or $\delta_A^-(r)$ with $A \in \text{next}(r)$

Theorem (Containment)

$$r \sqsubseteq s \Leftrightarrow (\nu(r) \Rightarrow \nu(s)) \wedge (\forall A \in \text{next}(r \sqsubseteq s)) \partial_A(r) \sqsubseteq \partial_A(s)$$

Notizen

Conclusion



- Generalize Brzozowski's derivative operator
- Extend Antimirov's algorithm for proving containment
- Provides a symbolic decision procedure that works with extended regular expressions on infinite alphabets
- Literals drawn from an effective boolean algebra
- Main contribution is to identify a finite set that covers all possibilities

Notizen

Regular Languages



The language $\llbracket r \rrbracket \subseteq \Sigma^*$ of a regular expression r is defined inductively by:

$$\begin{aligned}\llbracket \epsilon \rrbracket &= \{\epsilon\} \\ \llbracket A \rrbracket &= \{a \mid a \in A\} \\ \llbracket r+s \rrbracket &= \llbracket r \rrbracket \cup \llbracket s \rrbracket \\ \llbracket r \cdot s \rrbracket &= \llbracket r \rrbracket \cdot \llbracket s \rrbracket \\ \llbracket r^* \rrbracket &= \llbracket r \rrbracket \cdot \llbracket r^* \rrbracket \\ \llbracket r \&s \rrbracket &= \llbracket r \rrbracket \cap \llbracket s \rrbracket \\ \llbracket !r \rrbracket &= \overline{\llbracket r \rrbracket}\end{aligned}$$

Notizen

Nullable



The *nullable* predicate $\nu(r)$ indicates whether $\llbracket r \rrbracket$ contains the empty word, that is, $\nu(r)$ iff $\epsilon \in \llbracket r \rrbracket$.

$$\begin{aligned}\nu(\epsilon) &= \text{true} \\ \nu(A) &= \text{false} \\ \nu(r+s) &= \nu(r) \vee \nu(s) \\ \nu(r \cdot s) &= \nu(r) \wedge \nu(s) \\ \nu(r^*) &= \text{true} \\ \nu(r \&s) &= \nu(r) \wedge \nu(s) \\ \nu(!r) &= \neg \nu(r)\end{aligned}$$

Notizen

Brzowski Derivatives



$\partial_a(r)$ computes a regular expression for the left quotient $a^{-1}\llbracket r \rrbracket$.

$$\begin{aligned}\partial_a(\epsilon) &= \emptyset \\ \partial_a(A) &= \begin{cases} \epsilon, & a \in A \\ \emptyset, & a \notin A \end{cases} \\ \partial_a(r+s) &= \partial_a(r) + \partial_a(s) \\ \partial_a(r \cdot s) &= \begin{cases} \partial_a(r) \cdot s + \partial_a(s), & \nu(r) \\ \partial_a(r) \cdot s, & \neg \nu(r) \end{cases} \\ \partial_a(r^*) &= \partial_a(r) \cdot r^* \\ \partial_a(r \&s) &= \partial_a(r) \&\partial_a(s) \\ \partial_a(!r) &= !\partial_a(r)\end{aligned}$$

Notizen

First Symbols



Let $\text{first}(r) := \{a \mid aw \in \llbracket r \rrbracket\}$ be the set of first symbols derivable from regular expression r .

$$\begin{aligned} \text{first}(\epsilon) &= \emptyset \\ \text{first}(A) &= A \\ \text{first}(r+s) &= \text{first}(r) \cup \text{first}(s) \\ \text{first}(r \cdot s) &= \begin{cases} \text{first}(r) \cup \text{first}(s), & \nu(r) \\ \text{first}(r), & \neg \nu(r) \end{cases} \\ \text{first}(r^*) &= \text{first}(r) \\ \text{first}(r \&s) &= \text{first}(r) \cap \text{first}(s) \\ \text{first}(!r) &= \Sigma \setminus \{a \in \text{first}(r) \mid \partial_a(r) \neq \Sigma^*\} \end{aligned}$$

Notizen

First Literals



Let $\text{first}(r) := \{a \mid aw \in \llbracket r \rrbracket\}$ be the set of first symbols derivable from regular expression r .

$$\begin{aligned} \text{literal}(\epsilon) &= \emptyset \\ \text{literal}(A) &= \{A\} \\ \text{literal}(r+s) &= \text{literal}(r) \cup \text{literal}(s) \\ \text{literal}(r \cdot s) &= \begin{cases} \text{literal}(r) \cup \text{literal}(s), & \nu(r) \\ \text{literal}(r), & \neg \nu(r) \end{cases} \\ \text{literal}(r^*) &= \text{literal}(r) \\ \text{literal}(r \&s) &= \text{literal}(r) \cap \text{literal}(s) \\ \text{literal}(!r) &= \Sigma \cap \overline{\{A \in \text{literal}(r) \mid \partial_A(r) = \Sigma^*\}} \end{aligned}$$

Notizen

Coverage



Lemma (Coverage)

For all a , u , and r it holds that:

$$u \in \llbracket \partial_a(r) \rrbracket \Leftrightarrow \exists A \in \text{next}(r) : a \in A \wedge u \in \llbracket \delta_A^+(r) \rrbracket \wedge u \in \llbracket \delta_A^-(r) \rrbracket$$

Notizen

Termination



Theorem (Finiteness)

Let R be a finite set of regular inequalities. Define

$$F(R) = R \cup \{\partial_A(r \sqsubseteq s) \mid r \sqsubseteq s \in R, A \in \text{next}(r \sqsubseteq s)\}$$

For each r and s , the set $\bigcup_{i \in \mathbb{N}} F^{(i)}(\{r \sqsubseteq s\})$ is finite.

Notizen

Decision Procedure for Containment



$$\begin{array}{c} \text{(DISPROVE)} \\ \frac{\nu(r) \quad \neg\nu(s)}{\Gamma \vdash r \sqsubseteq s : \text{false}} \end{array} \quad \begin{array}{c} \text{(CYCLE)} \\ \frac{r \sqsubseteq s \in \Gamma}{\Gamma \vdash r \sqsubseteq s : \text{true}} \end{array}$$

$$\begin{array}{c} \text{(UNFOLD-TRUE)} \\ \frac{r \sqsubseteq s \notin \Gamma \quad \nu(r) \Rightarrow \nu(s) \quad \forall A \in \text{next}(r \sqsubseteq s) : \Gamma \cup \{r \sqsubseteq s\} \vdash \partial_A(r) \sqsubseteq \partial_A(s) : \text{true}}{\Gamma \vdash r \sqsubseteq s : \text{true}} \end{array}$$

$$\begin{array}{c} \text{(UNFOLD-FALSE)} \\ \frac{r \sqsubseteq s \notin \Gamma \quad \nu(r) \Rightarrow \nu(s) \quad \exists A \in \text{next}(r \sqsubseteq s) : \Gamma \cup \{r \sqsubseteq s\} \vdash \partial_A(r) \not\sqsubseteq \partial_A(s) : \text{false}}{\Gamma \vdash r \sqsubseteq s : \text{false}} \end{array}$$

Notizen

Prove and Disprove Axioms



$$\begin{array}{c} \text{(PROVE-IDENTITY)} \\ \Gamma \vdash r \sqsubseteq r : \text{true} \end{array} \quad \begin{array}{c} \text{(PROVE-EMPTY)} \\ \Gamma \vdash \emptyset \sqsubseteq s : \text{true} \end{array}$$

$$\begin{array}{c} \text{(PROVE-NULLABLE)} \\ \frac{\nu(s)}{\Gamma \vdash \epsilon \sqsubseteq s : \text{true}} \end{array} \quad \begin{array}{c} \text{(DISPROVE-EMPTY)} \\ \frac{\exists A \in \text{next}(r) : A \neq \emptyset}{\Gamma \vdash r \sqsubseteq \emptyset : \text{false}} \end{array}$$

Notizen

Soundness



Theorem (Soundness)

For all regular expression r and s :

$$\emptyset \vdash r \dot{\sqsubseteq} s : \top \Leftrightarrow r \sqsubseteq s$$

Notizen

Negative Derivatives



Counterexample

Let $r = \{a, b, c, d\} \cdot d^*$, $s = \{a, b, c\} \cdot d^* + \{b, c, d\} \cdot d^*$, and $A = \{a, b, c, d\}$, then

$$\delta_A^-(r) \dot{\sqsubseteq} \delta_A^+(s) \quad (4)$$

$$\delta_A^-(\{a, b, c, d\} \cdot d^*) \dot{\sqsubseteq} \delta_A^-(\{a, b, c\} \cdot d^*) + \delta_A^-(\{b, c, d\} \cdot d^*) \quad (5)$$

$$d^* \dot{\sqsubseteq} \emptyset + \emptyset \quad (6)$$

Notizen

Next Literals of an Inequality



Example

Let $r = \{a, b, c, d\} \cdot d^*$, $s = \{a, b, c\} \cdot c^* + \{b, c, d\} \cdot d^*$ then

$$\begin{aligned} \text{next}(r \dot{\sqsubseteq} s) &= \text{next}(\{a, b, c, d\} \cdot d^*) \times \text{next}(\{a, b, c\} \cdot c^* + \{b, c, d\} \cdot d^*) \\ &= \{\{a\}, \{b, c\}, \{d\}\} \end{aligned}$$

Notizen

Incomplete Containment



Conjecture

$$r \sqsubseteq s \Leftarrow (\nu(r) \Rightarrow \nu(s)) \wedge (\forall \mathbf{A} \in \text{litera}(r)) \delta_{\mathbf{A}}^+(r) \sqsubseteq \delta_{\mathbf{A}}^-(s)$$

Notizen

Notizen

Notizen
