

The Word-based Regular Expressions (originally CHeuSoV :-))

Aleksey Cheusov
vle@gmx.net

Minsk LUG meeting, 29 Dec. 2012, Belarus

Plan for my presentation

- Mathematics: RL definition and theorem
- Real World: applications of RL and some notes
- Linguistics: POS tagset, Semantic tagset, Text tagging, Tagged text corpus
- Wanted: (mumbling...)
- Mathematics: ExpRL definition, WRL definition
- Linguistics + Real World: WRE syntax and examples
- Вброс, драка, разбор завалов, развоз трупов

Definition: given a finite non-empty set of elements Σ (alphabet):

- \emptyset is a regular language.
- For any element $v \in \Sigma$, $\{v\}$ is a regular language.
- If A and B are regular languages, so is $A \cup B$.
- If A and B are regular languages, so is $\{ab | a \in A, b \in B\}$, where ab means string concatenation.
- If A is a regular language, so is A^* where $A^* = \emptyset \cup \{a_1 a_2 \dots a_n | a_i \in A, n > 0\}$.
- No other languages over Σ are regular.

Example:

- Let $\Sigma = \{a, b, c\}$. Then since aab and cc are members of Σ^* , $\{aab\}$ and $\{cc\}$ are regular languages. So is the union of these two sets $\{aab, cc\}$, and so is the concatenation of the two $\{aabcc\}$. Likewise, $\{aab\}^*$, $\{cc\}^*$ and $\{aab, cc, aabcc\}^*$ are regular languages.

Provable fact:

- Regular languages are closed under intersection and negation operations.

Definition: 5-tuple $(\Sigma, S, S_0, F, \delta)$ is a finite state automaton

- Σ is a finite non-empty set of elements (alphabet).
- S is a finite non-empty set of states.
- $S_0 \subseteq S$ is a set of *start* states.
- $F \subseteq S$ is a set of *final* states.
- $\delta : S \times \Sigma \rightarrow 2^S$ is a state transition function.

Theorem: \forall regular language R , \exists FSA f , such that $L(f) = R$; \forall FSA f , $L(f)$ is a regular language (L — language of FSA, that is a set of accepted inputs).

Real world. Regular expressions (hello UNIX!).

Regular language is a foundation for so called *regular expressions*, in most cases **alphabet** Σ is a set of characters (ASCII, Unicode etc.):

- POSIX basic regular expressions (*grep, sed, vi etc.*) and extended regular expressions (*grep -E, awk etc.*)
- Google re2, Yandex PIRE
- Perl, pcre, Ruby, Python (superset of regular language, and... extremely inefficient ;-))
- lex
- ...

Widely used extensions:

- submatch (depending on implementation may still be FSM but not FSA)
- backreferences (incompatible with regular language and FSM at all)

Penn part-of-speech tag set:

- **NN** noun, singular or mass (*apple, computer, fruit* etc.)
- **NNS** noun plural (*apples, computers, fruits* etc.)
- **CC** coordinating conjunction (*and, or*)
- **VB** verb, base form (*give, book, destroy* etc.)
- **VBD** verb, past tense form (*gave, booked, destroyed* etc.)
- **VBN** verb, past participle form (*given, booked, destroyed* etc.)
- **VBG** verb, gerund/present participle (*giving, booking, destroying* etc.)
- etc.

Semantic tag set:

- **LinkVerb** a verb that connects the subject to the complement(*seem, feel, look* etc.)
- **AnimateNoun** (*brother, son* etc.)
- etc.

Tagged sentence, POS tagging, semantic tagging

Examples:

- The book is red →

The_DT book_NN is_VBZ red_JJ →

The_DT book_NN/Object is_VBZ red_JJ/Color

Note: Words *book* and *red* are ambiguous.

- My son goes to school →

My_PRP\$ son_NN goes_VBZ to_IN school_NN →

My_PRP\$ son_NN/Person goes_VBZ to_TO
school_NN/Establishment

Questions: How to match tagged sentence or find portions of it?

Can regular expressions help? Do we need
regcomp(3)/regex(3)-like functionality?

Definition: given a non-empty set of elements Σ (alphabet) and a set of one-place predicates $P = \{P_1, P_2, \dots, P_k\}$,
 $P_i : \Sigma \rightarrow \{true, false\}$:

- \emptyset is an expanded regular language.
- For any element $v \in \Sigma$, $\{v\}$ is an expanded regular language.
- For any i $\{v | P_i(v) = true\}$ is an expanded regular language.
- Σ is an expanded regular language.
- If A and B are expanded regular languages, so is $A \cup B$.
- If A and B are expanded regular languages, so is $A \setminus B$.
- If A and B are expanded regular languages, so is $\{ab | a \in A, b \in B\}$, where ab means string concatenation.
- If A is a regular language, so is A^* where $A^* = \emptyset \cup \{a_1 a_2 \dots a_n | a_i \in A, n > 0\}$.
- No other languages over Σ are expanded regular.

Provable facts:

- Expanded regular languages are closed under intersection and negation operations.
- \forall expanded regular language R we can build a regular language R^* over alphabet Σ^* , such that $\exists f : \Sigma^* \rightarrow 2^\Sigma$ and $L(R) = L(R^*)$ (expanding all elements in $L(R^*)$ with a help of f). Thus, we can build expanded regular expression engine based on traditional FSA/FSM-based algorithms.

Definition: given

- W — set of words (character sequences), e.g.
"the", "apple", "123", ";", " C_2H_5OH ", ...
- T_{POS} — finite non-empty set of part-of-speech tags, e.g.
 $\{DT, NN, NNS, VBP, VBZ, \dots\}$
- T_{Sem} — finite non-empty set of semantic tags, e.g.
 $\{LinkVerb, Person, Object, TransitiveVerb, \dots\}$
- $D_{POS} : W \rightarrow 2^{T_{POS}}$ — POS dictionary, e.g.
 $D_{POS}(\text{"the"}) = \{DT\}$, $D_{POS}(\text{"book"}) = \{NN, VB, VBP\}$,
 $D_{POS}(\text{"and"}) = \{CC\}$
- $D_{Sem} : W \rightarrow 2^{T_{Sem}}$ — semantic dictionary, e.g.
 $D_{Sem}(\text{"son"}) = \{Person\}$,
 $D_{Sem}(\text{"mouse"}) = \{Animal, ComputerDevice\}$
- $EREs$ — finite set of POSIX extended regular expressions, e.g.
 $\{".*ing", ".multi.*al", "[A-Z][a-z]", "[0-9]^+" \dots\}$ etc.

(to be continued)

(continuation) the **word-based regular language**

$(T_{POS}, T_{Sem}, D_{POS}, D_{Sem}, EREs)$ is an **expanded regular language** over alphabet $W \times T_{POS} \times 2^{T_{Sem}}$ and one-place

predicates $P = \{P_{t_{POS}}^{check}, P_{t_{Sem}}^{check}, P_{t_{POS}}^{tagging}, P_{t_{Sem}}^{tagging}, P_{re}^{word}\}$, where

- $P_{t_{POS}}^{check}(w, \cdot, \cdot) = true$ if $t_{POS} \in D_{POS}(w)$, and *false* otherwise
- $P_{t_{Sem}}^{check}(w, \cdot, \cdot) = true$ if $t_{Sem} \in D_{Sem}(w)$, and *false* otherwise
- $P_{t_{POS}}^{tagging}(\cdot, tag^{POS}, \cdot) = true$ if $t_{POS} = tag^{POS}$, and *false* otherwise
- $P_{t_{Sem}}^{tagging}(\cdot, \cdot, tags^{Sem}) = true$ if $t_{Sem} \in tags^{Sem}$, and *false* otherwise
- $P_{re}^{word}(w, \cdot, \cdot) = true$ if POSIX extended regular expression $re \in EREs$ matches w , and *false* otherwise

The Word-based regular expressions (Finally!).

Syntax:

- **"word"** — word itself (P_{re}^{word}), e.g. "the", "2012-12-29" etc.
- **'regexp'** — words matched by specified regexp (P_{re}^{word})
- **Tag** — words tagged as tag_{POS} ($P_{t_{POS}}^{tagging}$), e.g. NN, DT, VB etc.
- **%Tag** — words tagged as tag_{Sem} ($P_{t_{Sem}}^{tagging}$), e.g. %Person, %Object, %LinkLerb etc.
- **_Tag** — words having as tag_{POS} in POS dictionary ($P_{t_{POS}}^{check}$), e.g. _NN, _DT, _VB etc.
- **@Tag** — words having as tag_{Sem} in semantic dictionary ($P_{t_{Sem}}^{check}$), e.g. @LinkVerb, @Object etc.
- **.** (dot) — any word with any POS and semantic tags
- **^** — beginning of the sentence
- **\$** — end of the sentence

(to be continued)

The Word-based regular expressions (Finally!).

Syntax (continuation):

- (R) — grouping like in mathematical expressions
- $\langle \text{num } R \rangle$ — submatch and extraction
- $R ?$ and $[R]$ — optional WRE
- $R ^*$ and $R ^+$ — possibly empty and non-empty repetitions
- $R \{n,m\}$, $R \{n,\}$ and $R \{,m\}$ — repetitions
- $R - S$ — subtraction
- $R \& S$ and R / S — intersection, $/$ is for single word WREs, $\&$ is for complex WREs
- $R | S$ — union
- $R S$ — concatenation
- $!R$ — negation ($L(!R)$ is equal to either $\Sigma \setminus L(R)$ or $\Sigma^* \setminus L(R)$ depending on a context of use)

(to be continued)

The Word-based regular expressions (Finally!).

Syntax (priorities from highest to lowest, continuation):

- , / and | in single word non-spaced WREs
- Prepositional unary operation !
- Postpositional unary operations {**n,m**}, '?', '+' and '*'
- (R) and <**num** R >
- R & S
- R – S
- R S
- R | S

(to be continued)

- How to select noun phrases (leftmost-longest match, only POS tags is our rules)

(DT | CD+)? RB * CC|JJ|JJR|JJS * (NN|NNS + | NP +)

Ex.: This absolutely stupid decision

Ex.: The best fuel cell

Ex.: Black and white colors

Ex.: Vasiliy Pupkin

- NER (Named Entity Recognition) for person names

$D_{Sem}(\text{"MrDr"}) = \{\text{"Mrs."}, \text{"Mr."}, \text{"Dr."}, \text{"Doctor"}, \text{"President"}, \text{"Dear"}, \dots\}$

@MrDr <1 '[A-Z][a-z]+'/!@MrDr_␣+>

Ex.: Doctor <1 Zhivago>

Ex.: Dr. <1 Vasiliy Pupkin>

Ex.: Mrs. <1 Kate>

but **not**

Ex.: Mr. <1 President>

- Question focus (object attributes)

\wedge "What" "is" "the" $\langle 1 \text{ @Attribute} \rangle$ "of" $\langle 2 .* \rangle$ "?"

Ex.: What is the $\langle 1 \text{ color} \rangle$ of $\langle 2 \text{ your book} \rangle$?

- Tiny definitions

$m4_define(\text{NounPhrase}, \text{"(see previous slide)"}) \wedge (\text{NounPhrase} \& (\dots)) \text{"is"} (\text{NounPhrase} \& (\dots))$

Ex.: What is the $\langle 1 \text{ color} \rangle$ of $\langle 2 \text{ your book} \rangle$?

- Sentiment analysis

$D_{Sem}(\text{"BadCharact"}) = \{\text{"sucks"}, \text{"stupid"}, \text{"crappy"}, \text{"shitty"}, \dots\}$

$D_{Sem}(\text{"GoodCharact"}) = \{\text{"rocks"}, \text{"awesome"}, \text{"excellent"}, \dots\}$

$D_{Sem}(\text{"OurProduct"}) = \{\text{"Linux"}, \text{"NetBSD"}, \text{"Ubuntu"}, \text{"AltLinux"}, \text{"iPad"}, \text{"Android"}, \dots\}$

$\langle 1 \text{ @BadCharact|@GoodCharact} \rangle \langle 2 \text{ @OurProduct} \rangle |$

$\langle 2 \text{ @OurProduct} \rangle \langle 1 \text{ @BadCharact|@GoodCharact} \rangle$

The Word-based Regular Expressions
is really cool DSL!