
Learning to Identify Regular Expressions that Describe Email Campaigns (Online Appendix)

Paul Prasse
Christoph Sawade
Niels Landwehr
Tobias Scheffer

PRASSE@CS.UNI-POTSDAM.DE
SAWADE@CS.UNI-POTSDAM.DE
LANDWEHR@CS.UNI-POTSDAM.DE
SCHEFFER@CS.UNI-POTSDAM.DE

University of Potsdam, Department of Computer Science, August-Bebel-Strasse 89, 14482 Potsdam, Germany

A. Definitions

Definition 2 (Regular Expressions). The set \mathcal{Y}_Σ of regular expressions over an ordered alphabet Σ is recursively defined as follows.

- Every $\mathbf{y}_j \in \Sigma \cup \{\epsilon, \cdot, \setminus S, \setminus e, \setminus w, \setminus d\}$, every range $\mathbf{y}_j = l_{min} \text{--} l_{max}$, where $l_{min}, l_{max} \in \Sigma$ and $l_{min} < l_{max}$, and their disjunction $[\mathbf{y}_1 \dots \mathbf{y}_k]$ are regular expressions.
- If $\mathbf{y}_1, \dots, \mathbf{y}_k \in \mathcal{Y}_\Sigma$ are regular expressions, so are the concatenation $\mathbf{y} = \mathbf{y}_1 \dots \mathbf{y}_k$, the disjunction $\mathbf{y} = \mathbf{y}_1 | \dots | \mathbf{y}_k$, $\mathbf{y} = \mathbf{y}_1^?$, $\mathbf{y} = (\mathbf{y}_1)$, and the repetitions $\mathbf{y} = \mathbf{y}_1^*$, $\mathbf{y} = \mathbf{y}_1^+$, $\mathbf{y} = \mathbf{y}_1\{l\}$, and $\mathbf{y} = \mathbf{y}_1\{l, u\}$, where $l, u \in \mathbb{N}$ and $l \leq u$.

We now define the syntax tree, the parse tree, and the matching lists for a regular expression \mathbf{y} and a string $x \in \Sigma^*$. The shorthand $(\mathbf{y} \rightarrow T_1, \dots, T_k)$ denotes the tree $T = (V, E, \Gamma, \leq)$ with root node $v_0 \in V$ labeled with $\Gamma(v_0) = \mathbf{y}$ and subtrees T_1, \dots, T_k . The order \leq maintains the subtree orderings \leq_i and defines the root node as the minimum over the set V and $v' \leq v''$ for all $v' \in V_i$ and $v'' \in V_j$, where $i < j$.

Definition 3 (Syntax Tree). The abstract syntax tree $T_{syn}^{\mathbf{y}}$ for a regular expression \mathbf{y} is recursively defined as follows. Let $T_{syn}^{\mathbf{y}_j} = (V_{syn}^{\mathbf{y}_j}, E_{syn}^{\mathbf{y}_j}, \Gamma_{syn}^{\mathbf{y}_j}, \leq_{syn}^{\mathbf{y}_j})$ be the syntax tree of the subexpression \mathbf{y}_j .

- If $\mathbf{y} \in \Sigma \cup \{\epsilon, \cdot, \setminus S, \setminus e, \setminus w, \setminus d\}$, or if $\mathbf{y} = l_{min} \text{--} l_{max}$, where $l_{min}, l_{max} \in \Sigma$, we define $T_{syn}^{\mathbf{y}} = (\mathbf{y} \rightarrow \emptyset)$.
- If $\mathbf{y} = (\mathbf{y}_1)$, where $\mathbf{y}_1 \in \mathcal{Y}_\Sigma$, we define $T_{syn}^{\mathbf{y}} = T_{syn}^{\mathbf{y}_1}$.

- If $\mathbf{y} = \mathbf{y}_1^*$, $\mathbf{y} = \mathbf{y}_1^+$, $\mathbf{y} = \mathbf{y}_1\{l, u\}$, or if $\mathbf{y} = \mathbf{y}_1\{l\}$, where $\mathbf{y}_1 \in \mathcal{Y}_\Sigma$, $l, u \in \mathbb{N}$, and there exist no $\mathbf{y}', \mathbf{y}'' \in \mathcal{Y}_\Sigma$ such that $\mathbf{y}_1 = \mathbf{y}' | \mathbf{y}''$ or $\mathbf{y}_1 = \mathbf{y}' \mathbf{y}''$, we define $T_{syn}^{\mathbf{y}} = (\mathbf{y} \rightarrow T_{syn}^{\mathbf{y}_1})$.
- If $\mathbf{y} = \mathbf{y}_1 \dots \mathbf{y}_k$, where $\mathbf{y}_j \in \mathcal{Y}_\Sigma$, and there exist no $\mathbf{y}', \mathbf{y}'' \in \mathcal{Y}_\Sigma$ such that $\mathbf{y}_j = \mathbf{y}' | \mathbf{y}''$ or $\mathbf{y}_j = \mathbf{y}' \mathbf{y}''$, we define $T_{syn}^{\mathbf{y}} = (\mathbf{y} \rightarrow T_{syn}^{\mathbf{y}_1}, \dots, T_{syn}^{\mathbf{y}_k})$.
- If $\mathbf{y} = \mathbf{y}_1 | \dots | \mathbf{y}_k$, where $\mathbf{y}_j \in \mathcal{Y}_\Sigma$, and there exist no $\mathbf{y}', \mathbf{y}'' \in \mathcal{Y}_\Sigma$ such that $\mathbf{y}_j = \mathbf{y}' | \mathbf{y}''$, or if $\mathbf{y} = [\mathbf{y}_1 \dots \mathbf{y}_k]$ and there exist no $\mathbf{y}', \mathbf{y}'' \in \mathcal{Y}_\Sigma$ such that $\mathbf{y}_j = \mathbf{y}' \mathbf{y}''$, we define $T_{syn}^{\mathbf{y}} = (\mathbf{y} \rightarrow T_{syn}^{\mathbf{y}_1}, \dots, T_{syn}^{\mathbf{y}_k})$.

Definition 4 (Parse Tree and Matching List). Given a syntax tree $T_{syn}^{\mathbf{y}} = (V_{syn}^{\mathbf{y}}, E_{syn}^{\mathbf{y}}, \Gamma_{syn}^{\mathbf{y}}, \leq_{syn}^{\mathbf{y}})$ of a regular expression \mathbf{y} with nodes $v \in V_{syn}^{\mathbf{y}}$ and a string $x \in L(\mathbf{y})$, a parse tree $T_{par}^{\mathbf{y}, x}$ and the matching lists $M^{\mathbf{y}, x}(v)$ for each $v \in V_{syn}^{\mathbf{y}}$ are recursively defined as follows. Let $T_{par}^{\mathbf{y}_j, x} = (V_{par}^{\mathbf{y}_j, x}, E_{par}^{\mathbf{y}_j, x}, \Gamma_{par}^{\mathbf{y}_j, x}, \leq_{par}^{\mathbf{y}_j, x})$ be the parse tree and $T_{syn}^{\mathbf{y}_j} = (V_{syn}^{\mathbf{y}_j}, E_{syn}^{\mathbf{y}_j}, \Gamma_{syn}^{\mathbf{y}_j}, \leq_{syn}^{\mathbf{y}_j})$ the syntax tree of the subexpression \mathbf{y}_j .

- If $\mathbf{y} = x$ and $x \in \Sigma \cup \{\epsilon\}$, we define $M^{\mathbf{y}, x}(v_0) = \{x\}$ and $T_{par}^{\mathbf{y}, x} = (\mathbf{y} \rightarrow \emptyset)$.
- If $\mathbf{y} = \cdot$ and $x \in \Sigma$, $\mathbf{y} = l_{min} \text{--} l_{max}$ and $l_{min} \leq x \leq l_{max}$, or if $\mathbf{y} \in \{\setminus S, \setminus w, \setminus e, \setminus d\}$ and x is either a non-whitespace character (everything but spaces, tabs, and line breaks), a word character (letters, digits, and underscores), a character in $\{., -, \#, +\}$ or a word character, or a digit,

respectively, we define

$$M^{\mathbf{y},x}(v) = \{x\} \text{ for all } v \in V_{syn}^{\mathbf{y}} \text{ and}$$

$$T_{par}^{\mathbf{y},x} = (\mathbf{y} \rightarrow T_{par}^{x,x}).$$

- If $\mathbf{y} = (\mathbf{y}_1)$ and $x \in \Sigma^*$, we define
 $M^{\mathbf{y},x}(v) = M^{\mathbf{y}_1,x}(v)$ for all $v \in V_{syn}^{\mathbf{y}}$ and
 $T_{par}^{\mathbf{y},x} = T_{par}^{\mathbf{y}_1,x}$

- If $\mathbf{y} = \mathbf{y}_1^*$, $x = x_1 \dots x_k$, and $k \geq 0$, or if
 $\mathbf{y} = \mathbf{y}_1^+$, and $k > 0$, or if
 $\mathbf{y} = \mathbf{y}_1\{l, u\}$, and $l \leq k \leq u$, or if
 $\mathbf{y} = \mathbf{y}_1\{l\}$, and $k = l$,
 where $x_i \in \Sigma^+$, and there exist no $\mathbf{y}', \mathbf{y}'' \in \mathcal{Y}_\Sigma$
 such that $\mathbf{y}_1 = \mathbf{y}'|\mathbf{y}''$ or $\mathbf{y}_1 = \mathbf{y}'\mathbf{y}''$, we define

$$M^{\mathbf{y},x}(v) = \begin{cases} \{x\} & , \text{ if } v = v_0 \\ \bigcup_{i=1}^k M^{\mathbf{y}_1, x_i}(v) & , \text{ if } v \in V_{syn}^{\mathbf{y}_1} \end{cases}, \text{ and}$$

$$T_{par}^{\mathbf{y},x} = (\mathbf{y} \rightarrow T_{par}^{\mathbf{y}_1, x_1}, \dots, T_{par}^{\mathbf{y}_1, x_k}).$$

- If $\mathbf{y} = \mathbf{y}_1 \dots \mathbf{y}_k$, $x = x_1 \dots x_k$,
 where $x_i \in \Sigma^*$, and there exist no $\mathbf{y}', \mathbf{y}'' \in \mathcal{Y}_\Sigma$
 such that $\mathbf{y}_j = \mathbf{y}'|\mathbf{y}''$ or $\mathbf{y}_j = \mathbf{y}'\mathbf{y}''$, we define

$$M^{\mathbf{y},x}(v) = \begin{cases} \{x\} & , \text{ if } v = v_0 \\ M^{\mathbf{y}_j, x_i}(v) & , \text{ if } v \in V_{syn}^{\mathbf{y}_j} \end{cases}, \text{ and}$$

$$T_{par}^{\mathbf{y},x} = (\mathbf{y} \rightarrow T_{par}^{\mathbf{y}_1, x_1}, \dots, T_{par}^{\mathbf{y}_k, x_k}).$$

- If $\mathbf{y} = \mathbf{y}_1 | \dots | \mathbf{y}_k$, $x \in \Sigma^*$
 and there exist no $\mathbf{y}', \mathbf{y}'' \in \mathcal{Y}_\Sigma$ such
 that $\mathbf{y}_j = \mathbf{y}' | \mathbf{y}''$, or if
 $\mathbf{y} = [\mathbf{y}_1 \dots \mathbf{y}_k]$, $x \in \Sigma^+$
 and there exist no $\mathbf{y}', \mathbf{y}'' \in \mathcal{Y}_\Sigma$ such
 that $\mathbf{y}_j = \mathbf{y}' \mathbf{y}''$, we define

$$M^{\mathbf{y},x}(v) = \begin{cases} \{x\} & , \text{ if } v = v_0 \\ M^{\mathbf{y}_j, x}(v) & , \text{ if } v \in V_{syn}^{\mathbf{y}_j}, \text{ and} \\ \emptyset & , \text{ otherwise} \end{cases}$$

$$T_{par}^{\mathbf{y},x} = (\mathbf{y} \rightarrow T_{par}^{\mathbf{y}_j, x}).$$

If $x \notin L(\mathbf{y})$, that is, no parse tree can be derived by the specification above, the empty sets $M^{\mathbf{y},x}(v) = \emptyset$ for all $v \in V_{syn}^{\mathbf{y}}$ and $T_{par}^{\mathbf{y},x} = \emptyset$ are returned. Otherwise, we denote the set of all parse trees and the unions of all matching lists for each $v \in V_{syn}^{\mathbf{y}}$ satisfying Definition 4 by $\mathcal{T}_{par}^{\mathbf{y},x}$ and $\mathcal{M}^{\mathbf{y},x}(v)$, respectively. Finally, the matching list $M^{\mathbf{y},\mathbf{x}}(v)$ for a set of strings \mathbf{x} for node $v \in V_{syn}^{\mathbf{y}}$ is defined as $M^{\mathbf{y},\mathbf{x}}(v) = \bigcup_{x \in \mathbf{x}} \mathcal{M}^{\mathbf{y},x}(v)$.

B. Used Features

Let M be a matching list and M_Σ be the set of characters in Σ which appear in the matching list M . The list of binary and continuous features, used to train *REx-SVM* is shown in Table 1.

Table 1. Important features used to train REx-SVM.

Feature	Description
$\llbracket \varepsilon \in M \rrbracket$	Matching list contains the empty string?
$\llbracket \forall \mathbf{x} \in M \mathbf{x} = 1 \rrbracket$	All elements of the matching list have the length one?
$\llbracket \exists i \in \mathbb{N} \forall \mathbf{x} \in M \mathbf{x} = i \rrbracket$	All elements of the matching list have the same length?
$\frac{26}{ M_\Sigma \cap \{A, \dots, Z\} }$	Portion of characters A-Z in the matching list
$\frac{26}{ M_\Sigma \cap \{a, \dots, z\} }$	Portion of characters a-z in the matching list
$\frac{10}{ M_\Sigma \cap \{0, \dots, 9\} }$	Portion of characters 0-9 in the matching list
$\frac{6}{ M_\Sigma \cap \{A, \dots, F\} }$	Portion of characters A-F in the matching list
$\frac{6}{ M_\Sigma \cap \{a, \dots, f\} }$	Portion of characters a-f in the matching list
$\frac{6}{ M_\Sigma \cap \{G, \dots, Z\} }$	Portion of characters G-Z in the matching list
$\frac{20}{ M_\Sigma \cap \{g, \dots, z\} }$	Portion of characters g-z in the matching list
$\llbracket \forall x \in M_\Sigma x \notin \{A, \dots, Z\} \rrbracket$	No characters of A-Z in the matching list?
$\llbracket \forall x \in M_\Sigma x \notin \{a, \dots, z\} \rrbracket$	No characters of a-z in the matching list?
$\llbracket \forall x \in M_\Sigma x \notin \{0, \dots, 9\} \rrbracket$	No characters of 0-9 in the matching list?
$\llbracket \forall x \in M_\Sigma x \notin \{a, \dots, f\} \rrbracket$	No characters of a-f in the matching list?
$\llbracket \forall x \in M_\Sigma x \notin \{A, \dots, F\} \rrbracket$	No characters of A-F in the matching list?
$\llbracket M_\Sigma \cap \{-, /, ?, =, ., @, : \} > 0 \rrbracket$	Matching list contains URL/Email characters?
$\llbracket \forall \mathbf{x} \in M \mathbf{x} \geq 1 \wedge \mathbf{x} \leq 5 \rrbracket$	Length of strings in the matching list is between 1 and 5?
$\llbracket \forall \mathbf{x} \in M \mathbf{x} \geq 6 \wedge \mathbf{x} \leq 10 \rrbracket$	Length of strings in the matching list is between 5 and 10?
$\llbracket \forall \mathbf{x} \in M \mathbf{x} \geq 11 \wedge \mathbf{x} \leq 20 \rrbracket$	Length of strings in the matching list is between 10 and 20?
$\llbracket \forall \mathbf{x} \in M \mathbf{x} > 20 \rrbracket$	Length of strings in the matching list is higher than 20?
$\llbracket M = 0 \rrbracket$	Matching list is empty?