

Theory and Practice of Data Cleaning

Introduction to Regular Expressions: From **Theory** to Practice



Introduction to Regular Expressions (Regex)

Theory & Practice

- **Theory** of regular expressions:
 - Brief introduction where regular expressions come from ...
- **Practice** of regular expressions:
 - What you need to know to get started with regex in practice!
- **Demonstration** of regular expressions

Why study **regular expressions**?

- Widely used **in practice**:
 - A bit like *wildcards* (e.g. find all csv file: `*.csv`)
 - ... but much **more powerful** (“*wildcards on steroids*”) !
- Used to *match, extract, find-and-replace* data, e.g.,
 - ... in text editors
 - ... scripting and programming (Bash, Python, Perl, R, ..., Java, ...)
 - ... screen scraping and other data extraction applications

Why study regular expressions for **data cleaning**?

- Useful to **match** (*assess*) and **transform** (*clean*) data:

- OpenRefine Expression Language (GREL)
- Use in scripting languages for **data cleaning**
- ... *workflow automation*

- Example: *ISO 8601 date format*:

- **YYYY-MM-DD**

- ... *vs other (common) date formats*:

- MM/DD/YY
- DD.MM.YYYY
- ...

PUBLIC SERVICE ANNOUNCEMENT:

OUR DIFFERENT WAYS OF WRITING DATES AS NUMBERS CAN LEAD TO ONLINE CONFUSION. THAT'S WHY IN 1988 ISO SET A GLOBAL STANDARD NUMERIC DATE FORMAT.

THIS IS **THE** CORRECT WAY TO WRITE NUMERIC DATES:

2013-02-27

THE FOLLOWING FORMATS ARE THEREFORE DISCOURAGED:

02/27/2013 02/27/13 27/02/2013 27/02/13
20130227 2013.02.27 27.02.13 27-02-13
27.2.13 2013.II.27. 27²/13 2013.158904109
MMXIII-II-XXVII MMXIII ^{LVII}CCCLXV 1330300800
((3+3)×(111+1)-1)×3/3-1/3³ 2013 2-27-13 miss
10/11011/1101 02/27/20/13 0 1 2 3 4 5 6 7 8

Introduction to Regular Expressions (Regex)

Theory & Practice

- **Theory** of regular expressions:
 - Brief introduction where regular expressions come from ...
- Practice of regular expressions:
 - What you need to know to get started with regex in practice!
- Demonstration of regular expressions

Theory of Regular Expressions

- **Regular expression** (regex): in *theoretical computer science* (esp. *formal language theory*):
 - A formal expression that defines a *search pattern*
 - ... used to *match* (or *recognize*) a strings

Theory of Regular Expressions

- **Formal definition:**
- Base elements:
 - \emptyset *empty set*, ϵ *empty string*, and Σ *alphabet* of characters
- Given regular expressions R and S , the following are also regular expressions:
 - $R \mid S$ *alternation*
 - RS *concatenation*
 - R^* *Kleene star*
 - (R) *parentheses* (can be omitted with *precedence rules*)

Regular Languages in the Chomsky Hierarchy

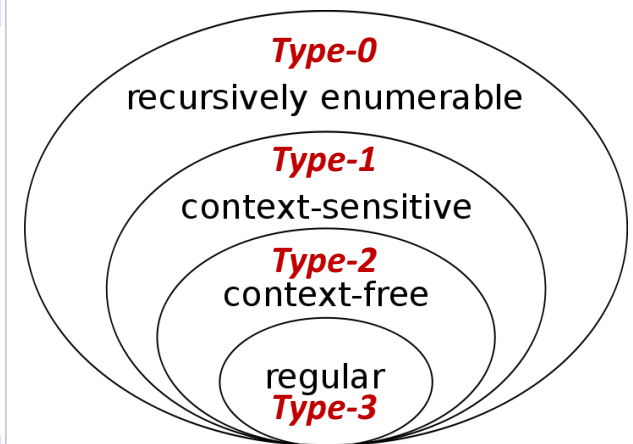


Automata theory: formal languages and formal grammars [hide]			
Chomsky hierarchy	Grammars	Languages	Abstract machines
Type-0	Unrestricted	Recursively enumerable	Turing machine
—	(no common name)	Decidable	Decider
Type-1	Context-sensitive	Context-sensitive	Linear-bounded
—	Positive range concatenation	Positive range concatenation*	PTIME Turing Machine
—	Indexed	Indexed*	Nested stack
—	—	—	Thread automaton
—	Linear context-free rewriting systems	Linear context-free rewriting language	restricted Tree stack automaton
—	Tree-adjoining	Tree-adjoining	Embedded pushdown
Type-2	Context-free	Context-free	Nondeterministic pushdown
—	Deterministic context-free	Deterministic context-free	Deterministic pushdown
—	Visibly pushdown	Visibly pushdown	Visibly pushdown
Type-3	Regular	Regular	Finite
—	—	Star-free	Counter-free (with aperiodic finite monoid)
—	Non-recursive	Finite	Acyclic finite

Each category of languages, except those marked by a *, is a proper subset of the category directly above it.
Any language in each category is generated by a grammar and by an automaton in the category in the same line.

Categories: Formal languages | Finite automata

The "regular" in regular expression



https://en.wikipedia.org/wiki/Regular_language

Regular Grammars

Example: **floating point numbers** such as $-0.314159265e+1$... can be **generated** by a **right regular grammar** G with $N = \{S, A, B, C, D, E, F\}$, $\Sigma = \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9, +, -, ., e\}$,



Production rules $P =$

$S \rightarrow +A$	$A \rightarrow 0A$	$B \rightarrow 0C$	$C \rightarrow 0C$	$D \rightarrow +E$	$E \rightarrow 0F$	$F \rightarrow 0F$
$S \rightarrow -A$	$A \rightarrow 1A$	$B \rightarrow 1C$	$C \rightarrow 1C$	$D \rightarrow -E$	$E \rightarrow 1F$	$F \rightarrow 1F$
$S \rightarrow A$	$A \rightarrow 2A$	$B \rightarrow 2C$	$C \rightarrow 2C$	$D \rightarrow E$	$E \rightarrow 2F$	$F \rightarrow 2F$
	$A \rightarrow 3A$	$B \rightarrow 3C$	$C \rightarrow 3C$		$E \rightarrow 3F$	$F \rightarrow 3F$
	$A \rightarrow 4A$	$B \rightarrow 4C$	$C \rightarrow 4C$		$E \rightarrow 4F$	$F \rightarrow 4F$
	$A \rightarrow 5A$	$B \rightarrow 5C$	$C \rightarrow 5C$		$E \rightarrow 5F$	$F \rightarrow 5F$
	$A \rightarrow 6A$	$B \rightarrow 6C$	$C \rightarrow 6C$		$E \rightarrow 6F$	$F \rightarrow 6F$
	$A \rightarrow 7A$	$B \rightarrow 7C$	$C \rightarrow 7C$		$E \rightarrow 7F$	$F \rightarrow 7F$
	$A \rightarrow 8A$	$B \rightarrow 8C$	$C \rightarrow 8C$		$E \rightarrow 8F$	$F \rightarrow 8F$
	$A \rightarrow 9A$	$B \rightarrow 9C$	$C \rightarrow 9C$		$E \rightarrow 9F$	$F \rightarrow 9F$
	$A \rightarrow .B$		$C \rightarrow eD$			$F \rightarrow \epsilon$
	$A \rightarrow B$		$C \rightarrow \epsilon$			

Regular Grammars

Example: **floating point numbers** such as **-0.314159265e+1**
 ... can be **generated** by a **right regular grammar G** with
 $N = \{S, A, B, C, D, E, F\}$, $\Sigma = \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9, +, -, ., e\}$,



Production rules P =

$S \rightarrow +A$	$A \rightarrow 0A$	$B \rightarrow 0C$	$C \rightarrow 0C$	$D \rightarrow +E$	$E \rightarrow 0F$	$F \rightarrow 0F$
$S \rightarrow -A$	$A \rightarrow 1A$	$B \rightarrow 1C$	$C \rightarrow 1C$	$D \rightarrow -E$	$E \rightarrow 1F$	$F \rightarrow 1F$
$S \rightarrow A$	$A \rightarrow 2A$	$B \rightarrow 2C$	$C \rightarrow 2C$	$D \rightarrow E$	$E \rightarrow 2F$	$F \rightarrow 2F$
	$A \rightarrow 3A$	$B \rightarrow 3C$	$C \rightarrow 3C$		$E \rightarrow 3F$	$F \rightarrow 3F$
	$A \rightarrow 4A$	$B \rightarrow 4C$	$C \rightarrow 4C$		$E \rightarrow 4F$	$F \rightarrow 4F$
	$A \rightarrow 5A$	$B \rightarrow 5C$	$C \rightarrow 5C$		$E \rightarrow 5F$	$F \rightarrow 5F$
	$A \rightarrow 6A$	$B \rightarrow 6C$	$C \rightarrow 6C$		$E \rightarrow 6F$	$F \rightarrow 6F$
	$A \rightarrow 7A$	$B \rightarrow 7C$	$C \rightarrow 7C$		$E \rightarrow 7F$	$F \rightarrow 7F$
	$A \rightarrow 8A$	$B \rightarrow 8C$	$C \rightarrow 8C$		$E \rightarrow 8F$	$F \rightarrow 8F$
	$A \rightarrow 9A$	$B \rightarrow 9C$	$C \rightarrow 9C$		$E \rightarrow 9F$	$F \rightarrow 9F$
	$A \rightarrow .B$		$C \rightarrow eD$			$F \rightarrow \varepsilon$
	$A \rightarrow B$		$C \rightarrow \varepsilon$			

- Not very handy in practice ...
- **Regular expressions** to the rescue!

$[-+]?[0-9]^*\backslash.?[0-9]+([eE][-+]?[0-9]+)?$