

# **CE4045 CZ4045 SC4002**

## **Natural Language Processing**

### **Regular Expressions**

Dr. Sun Aixin



# Regular Expressions (RE)

- RE is a language for specifying text search strings.
  - Used in every computer language, word processor, and text processing tools
  - Formally, a RE is an **algebraic notation** for characterizing a **set** of strings.
- Particularly useful for searching in texts, when we have a **pattern** to search for and a **corpus** of texts to search through.
  - A RE search function will either return **the first match**, or search through the corpus and return **all texts that match the pattern**
  - The corpus can be a **single document** or a **collection**
  - Example: search for **postcode** in **5 documents**

The screenshot shows a web-based Regular Expression testing tool. The 'Expression' field contains the regex `/([A-Z])\w+/'g`. The 'Text' field contains a paragraph of text about RegExr. The 'Tests' tab shows 29 matches. The 'Tools' section includes tabs for 'Replace', 'List', 'Details', and 'Explain'. A tooltip explains the capturing group #1 and the character set [A-Z].

# Basic Regular Expression Patterns

- The simplest RE is a **sequence of simple characters**, like `/test/`
  - Putting characters in sequence is called **concatenation**.
  - `//` is not part of the RE, but the notation to indicate RE.
  - `/test/` will match test in text, but not match Test, or TEST.
  - Regular expressions are **case sensitive**.



RE	Example patterns matched
<code>/woodchucks/</code>	“interesting links to <u>woodchucks</u> and lemurs”
<code>/a/</code>	“M <u>a</u> ry Ann stopped by Mona’s”
<code>/!/</code>	“You’ve left the burglar behind again <u>!</u> ” said Nori

- In our examples, we generally underline the **first exact part** of the match.
  - If the word woodchucks appears to be the first word in a sentence, then *Woodchucks* will not be match.

# Disjunction of characters to match

- The string of characters inside the square braces [ ] specifies a disjunction of characters to match
  - Square brackets [ ] matches **any single character** from within the class

RE	Match	Example Patterns
<code>/[wW]oodchuck/</code>	Woodchuck or woodchuck	“ <u>Woodchuck</u> ”
<code>/[abc]/</code>	‘a’, ‘b’, or ‘c’	“In uomini, in soldati”
<code>/[1234567890]/</code>	any single digit	“plenty of <u>7</u> to 5”

- If there is a well-defined sequence associated with a set of characters, [ ] can be used with the dash (-) to specify any one character in a range.
  - `/[A-Z]/` matches an upper-case letter
  - `/[a-z]/` matches a lower-case letter
  - `/[0-9]/` matches a single digit



# More on square braces in RE

- The square braces can also be used to specify what a single character cannot be matched with caret **^**.
- If the caret **^** is **the first symbol after the open square brace** **[**, the resulting pattern is negated.
  - Pattern **/[^a]/** matches any single character (including special characters) except **a**.
  - This is only true when the caret is the first symbol after the open square brace. If it occurs anywhere else, it usually stands for a caret

RE	Match (single characters)	Example Patterns Matched
<b>/[^A-Z]/</b>	not an upper-case letter	“O <b>y</b> fn pripetchik”
<b>/[^Ss]/</b>	neither ‘S’ nor ‘s’	“ <b>I</b> have no exquisite reason for’t”
<b>/[^.]/</b>	not a period	“ <b>o</b> ur resident Djinn”
<b>/[e^]/</b>	either ‘e’ or ‘^’	“look up <b>^</b> now”
<b>/a^b/</b>	the pattern ‘a^b’	“look up <b>a^b</b> now”

# Counters in RE, and wildcard expression

- Question mark **?** matches zero or one appearance of the preceding item
- Kleene **\*** (generally pronounced “cleany star”) matches zero or more occurrences of the immediately previous character or regular expression
- Kleene **+** matches one or more occurrences of the immediately preceding character or regular expression.
- Examples
  - **/colou?r/** matches color or colour
  - **/ba\*!/** matches b!, ba!, baaa!, baaaaa! and more such strings
  - **/ba+!/** matches ba!, baa!, baaaaa! and more such strings
  - **/[ab]+/** means “one or more *a*’s or *b*’s” like aaaa, or ababbb or bbbb
- The period (**/./**) is a wildcard expression that matches any single character (except a carriage return), e.g., **/beg.n/** matches begin, beg’n, begun
  - The wildcard used together with the Kleene star **/.\*/** means “any string of characters”.



# Anchors in RE

- **Anchors** are special characters that anchor regular expressions to particular places in a string
- The caret **^** matches the **start of a line**.
  - The pattern **/^The/** matches the word **The** only at the start of a line.
- The dollar sign **\$** matches the **end of a line**.
  - **/^The dog\.\$/** matches a line that contains only the phrase “**The dog.**”
  - Backslash dot **\.** is escaped character, means **the full stop or period, not the wildcard**.
- Two other anchors: word boundary **\b** and non-word boundary **\B**
  - A “word” for the purposes of a regular expression is defined as any sequence of **digits**, **underscores**, or **letters**, based on the definition of “words” in programming languages.
  - **/\bthe\b/** matches the word **the** but not the word **other**
  - **/\b99\b/** will match “There are **99** bottles” (because 99 follows a space) but not “There are **299** bottles” (since 99 follows a number). But it will match **\$99**, since 99 follows a dollar sign (not part of the word definition above).



# Disjunction, Grouping, and Precedence

➤ The disjunction operator, also called the pipe symbol `|`, means “or”.

- The pattern `/cat|dog/` matches **either** the string cat **or** the string dog

➤ Parenthesis operators ( and ).

- Enclosing a pattern in parentheses makes it **act like a single character** for the purposes of neighboring operators like the pipe `|` and the Kleene\*.
- `/gupp(y|ies)/` would specify that we meant the disjunction only to apply to the suffixes **y** and **ies**. The RE matches guppy or guppies
- `/(abb)+/` matches abb, abbabb, abbabbabb ....

➤ Operator precedence

- Parenthesis ( )
- Counters `* + ? { }`
- Sequences and anchors the `^my end$`
- Disjunction `|`

REs always match the largest string they can!





# More operators

RE	Expansion	Match	First Matches
\d	[0-9]	any digit	Party_of_5
\D	[^0-9]	any non-digit	Blue_moon
\w	[a-zA-Z0-9_]	any alphanumeric/underscore	Daiyu
\W	[^\w]	a non-alphanumeric	!!!!
\s	[\r\t\n\f]	whitespace (space, tab)	
\S	[^\s]	Non-whitespace	in_Concord

RE	Match
*	zero or more occurrences of the previous char or expression
+	one or more occurrences of the previous char or expression
?	exactly zero or one occurrence of the previous char or expression
{n}	<i>n</i> occurrences of the previous char or expression
{n,m}	from <i>n</i> to <i>m</i> occurrences of the previous char or expression
{n,}	at least <i>n</i> occurrences of the previous char or expression
{,m}	up to <i>m</i> occurrences of the previous char or expression

RE	Match	First Patterns Matched
\*	an asterisk “*”	“K_A*P*L*A*N”
\.	a period “.”	“Dr. Livingston, I presume”
\?	a question mark	“Why don’t they come and lend a hand?”
\n	a newline	
\t	a tab	

# A working example

- Suppose we wanted to write a RE to find the English article *the* in text.

Expression

`/the/g` Missing “The”, and matching “there”

Text Tests

The word the can be a substring of "there" or "the25", or even "the\_example".

Expression

`/[tT]he/g` Matching both “the” and “The”, but still matching “there”

Text Tests

The word the can be a substring of "there" or "the25", or even "the\_example".



# A working example

- Suppose we wanted to write a RE to find the English article *the* in text.

Expression

```
/\b[tT]he\b/g
```

Added word boundary, missing two other cases

Text Tests

The word the can be a substring of "there" or "the25", or even "the\_example".

Expression

```
/^[^a-zA-Z][tT]he[^a-zA-Z]/g
```

No other alphabetic letters on either side of "the"

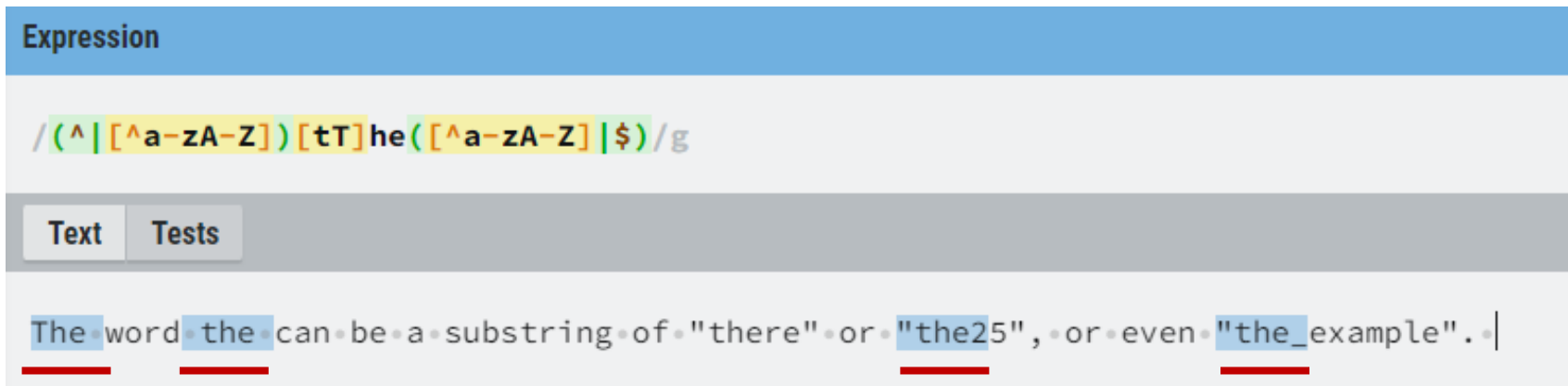
Text Tests

The word the can be a substring of "there" or "the25", or even "the\_example".



# A working example

➤ Suppose we wanted to write a RE to find the English article *the* in text.



- A more complicated (but not perfect) version
- Match “the” or “The” either in text or at the beginning of the line, or at the end of line.
- There is no alphabetic letters on either side of the word “the”.
- The final version can be application dependent.

# False Positive and False Negative

---

- The whole process is to fix two kinds of errors:
  - **False positives:** strings that we incorrectly matched like **other** or **there**,
  - **False negatives:** strings that we incorrectly missed, like **The**.
- Addressing these two kinds of errors is common in implementing language processing systems.
- Reducing the overall error rate for an application thus involves
  - Increasing **precision** (minimizing false positives)
  - Increasing **recall** (minimizing false negatives)



# Summary

---

- RegExr: an online tool to learn, build, & test Regular Expressions
  - <http://regexr.com/>
- Java RegEx API and Tutorial
  - <http://docs.oracle.com/javase/8/docs/api/java/util/regex/package-summary.html>
  - <http://docs.oracle.com/javase/tutorial/essential/regex/>
- Reference: <https://web.stanford.edu/~jurafsky/slp3/>
  - **Chapter 2**, Regular Expressions, Text Normalization, Edit Distance

