

Natural Language Processing

Session 2: Basic string processing and regular expressions

Sascha Göbel
goebel@hertie-school.org

September 17, 2025

Basic string processing

Strings

- Record textual information
- Collections of one-character strings
- Sequences that maintain a left-to-right order
- In NL elements of these sequences are not independent

Basic string processing

Strings

- Record textual information
- Collections of one-character strings
- Sequences that maintain a left-to-right order
- In NL elements of these sequences are not independent

In Python

- Enclosed in single or double quotes
- Immutable (cannot be changed in place)
- Modification, extraction, splitting operations via built-in methods and functions

Regular Expressions (regex)

- A set of characters to match and find substrings in a text
- A tool for describing text patterns

Why regular expressions?

- Pattern-based/rule-based approach to NLP
- Generalized modification, extraction, splitting operations
- Useful in the absence of large amounts of data
- Useful to encode domain knowledge

In Python

- via `re` library
- using `re.match()`, `re.search()`, `re.sub()`, `re.findall()`, `re.split()`

Method	Purpose
<code>re.compile(<i>pattern</i>)</code>	Compile regex pattern
<code>re.match(<i>pattern</i>, <i>string</i>)</code>	Find first match from beginning of the string
<code>re.search(<i>pattern</i>, <i>string</i>)</code>	Find first match throughout the string
<code>re.findall(<i>pattern</i>, <i>string</i>)</code>	Find all matches throughout the string
<code>re.sub(<i>pattern</i>, <i>repl</i>, <i>string</i>)</code>	Replace all matches throughout the string
<code>re.split(<i>pattern</i>, <i>string</i>)</code>	Split string by all pattern occurrences

Regex patterns and grammar

Anchors

- Match a position before, after, or between characters

Pattern	Match
<code>^</code>	beginning of string
<code>\$</code>	end of string
<code>\b</code>	word boundary

Character classes

- Match only one out of several characters

Pattern	Match
<code>[A-Z]</code>	upper case letters
<code>[a-z]</code>	lower case letters
<code>[0-9]</code>	digits
<code>[^A-Z]</code>	characters that are not upper case letters

Aliases

- Shorthand character classes

Pattern	Match
<code>\d</code>	any digit, e.g., <code>[0-9]</code>
<code>\D</code>	any non digit, e.g., <code>[^0-9]</code>
<code>\w</code>	any alphanumeric, e.g., <code>[A-Za-z0-9_]</code>
<code>\W</code>	any non-alphanumeric, e.g., <code>[^A-Za-z0-9_]</code>
<code>\s</code>	any whitespace, e.g., <code>[\t\r\n\f]</code>
<code>\S</code>	any non-whitespace, e.g., <code>[^\t\r\n\f]</code>

Quantifiers

- Repetition operators

Pattern	Match
<code>?</code>	zero or one instance of the previous character
<code>*</code>	zero or more occurrences of the previous character
<code>+</code>	one or more occurrences of the previous character
<code>{n}</code>	exactly n occurrences of the previous character
<code>{n,m}</code>	from n to m occurrences of the previous character
<code>{n,}; {,m}</code>	at least n; at least m occurrences of the previous character

Grouping and capturing

- Group parts of a regex together and potentially backreference

Pattern	Match
<i>(pattern)</i>	group <i>pattern</i> and capture, automatically numbered
<i>\1</i>	Recall/reference the captured group by number
<i>(?:pattern)</i>	group <i>pattern</i> without capturing

Lookaround

- Match conditional on context

Pattern	Match
<i>match(?=context)</i>	positive lookahead
<i>match(?!context)</i>	negative lookahead
<i>(?<=context)match</i>	positive lookbehind
<i>(?<!=context)match</i>	negative lookbehind

Special characters

- Characters with special meaning, escaped with `\` for literal use

Characters	Meaning
.	wildcard, matches (almost) any character
?, *, +, { }	quantifiers
	disjunction
()	groupings
[]	character classes
\	escape sequences, aliases, backreferences, non-printable characters

Precedence

- Order of evaluation

Order	Operator
1.	\ (escaped characters)
2.	[] (character classes)
3.	() (groups)
4.	?, *, +, {} (quantifiers)
5.	literal characters, ^, \$, \b
6.	(alternation)