Chapter NLP:I

- I. Natural Language Processing Basics
 - □ String Processing
 - □ Grammars
 - □ Regular Expressions

String Processing

Overview

String processing forms the basis of many tasks, algorithms, and evaluation methods in NLP.

String problems:

- □ Sorting
- Manipulation

Exact / inexact matching Search and indexing, similarity and distance

Alignment

Longest common subsequence

Parsing Splitting

Compression

Data structures:

Buffer

- Inverted index
- □ Trie, suffix tree, suffix array

Computational models:

- □ Finite-state machines
- Dynamic programming

□ String processing combines methods from formal language theory, compilers, and bioinformatics.

Terminology

\Box Alphabet Σ .

An alphabet Σ is a finite non-empty set of letters or symbols.

\Box Word w.

A word w is a finite sequence of symbols from Σ . The length of a word |w| is the number of its symbols.

 ε denotes the empty word; it is the only word with length 0.

 Σ^* denotes the set of all words over Σ .

□ Language *L*.

A language L is a set of words over an alphabet Σ .

\Box Grammar G.

A grammar G is a calculus for defining a language – that is, a set of rules that can be used to derive words. The language belonging to G consists of all derivable, terminal words.

Terminology

D Alphabet Σ .

An alphabet Σ is a finite non-empty set of letters or symbols.

\Box Word w.

A word w is a finite sequence of symbols from Σ . The length of a word |w| is the number of its symbols.

 ε denotes the empty word; it is the only word with length 0.

 Σ^* denotes the set of all words over $\Sigma.$

□ Language *L*.

A language *L* is a set of words over an alphabet Σ .

\Box Grammar G.

A grammar *G* is a calculus for defining a language – that is, a set of rules that can be used to derive words. The language belonging to *G* consists of all derivable, terminal words.

- When defining language characteristics, a distinction is made between different levels of abstraction. Level 1 deals with the notation of symbols, while Level 2 deals with the syntactic structure of the language.
 [Exkurs: Programmiersprachen]
- □ To distinguish between the different levels of grammar usage, the following terms can be used:
 - Level 1: Alphabet, symbol, word, language
 - Level 2: Vocabulary, token, sentence, language
- □ The words {alphabet, vocabulary}, {symbol, token} or {word, sentence} are the respective equivalents of the elementary symbolic level and the syntactic level.

Definition 1 (Grammar)

A grammar is a quadruple $G = (N, \Sigma, P, S)$ with

- N =finite set of non-terminal symbols
- Σ = finite set of terminal symbols, $N \cap \Sigma = \emptyset$
- P = finite set of productions or rules

 $P \subset (N \cup \Sigma)^* N (N \cup \Sigma)^* \times (N \cup \Sigma)^*$

S =start symbol, $S \in N$

Definition 1 (Grammar)

A grammar is a quadruple $G = (N, \Sigma, P, S)$ with

- N =finite set of non-terminal symbols
- Σ = finite set of terminal symbols, $N \cap \Sigma = \emptyset$
- P = finite set of productions or rules

$$P \subset \underbrace{(N \cup \Sigma)^* N (N \cup \Sigma)^*}_{A} \times \underbrace{(N \cup \Sigma)^*}_{Ab}$$

S =start symbol, $S \in N$

- A rule consists of a left side (premise) and a right side (conclusion), each of which is a word consisting of terminals and non-terminals. The left side must contain at least one non-terminal, and unlike the left side, the right side can also be the empty word. [Wikipedia]
- □ A rule can be applied to a word consisting of terminals and non-terminals, whereby any occurrence of the left side of the rule in the word is replaced by the right side of the rule: $w \to w'$.
- Given the rule $w \to w'$, then w, w' are in the so-called *transitive relation* \to_G . A sequence of rule applications is called a *derivation*.

Definition 2 (Generated Language)

The language L(G) generated by a grammar $G = (N, \Sigma, P, S)$ contains exactly those words that consist only of terminal symbols and can be derived from the start symbol with a finite number of steps:

$$L(G) := \{ w \in \Sigma^* \mid S \to_G^* w \}$$

 \rightarrow^*_G denotes the arbitrary application of the productions in *G*, i.e., the reflexive-transitive hull of the transition relation \rightarrow_G .

Definition 2 (Generated Language)

The language L(G) generated by a grammar $G = (N, \Sigma, P, S)$ contains exactly those words that consist only of terminal symbols and can be derived from the start symbol with a finite number of steps:

$$L(G) := \{ w \in \Sigma^* \mid S \to_G^* w \}$$

 \rightarrow^*_G denotes the arbitrary application of the productions in *G*, i.e., the reflexive-transitive hull of the transition relation \rightarrow_G .

Example:

 $G = (N, \Sigma, P, S)$ mit $N = \{S, A, B\}$, $\Sigma = \{a, b\}$ und folgenden Produktionen:

S	\rightarrow	ABS	BA	\rightarrow	AB	Ab	\rightarrow	ab
S	\rightarrow	ε	BS	\rightarrow	b	Aa	\rightarrow	aa
			Bb	\rightarrow	bb			

- It is conventional to use uppercase letters to denote non-terminal symbols and lowercase letters to denote terminal symbols.
- □ Another grammar that generates the same language as in the example is: $N = \{S, A, B\}, \Sigma = \{a, b\}, P = \{S \rightarrow aSb, S \rightarrow \varepsilon\}$

Chomsky Hiearchy

Grammars are divided into four classes based on the complexity of the languages they generate.

- $\hfill\square$ Type 0 \sim recursively enumerable.
- $\hfill\square$ Type 1 \sim context-sensitive.
- $\hfill\square$ Type 2 \sim context-free.
- $\hfill\square$ Type 3 \sim regular.

Chomsky Hiearchy

Grammars are divided into four classes based on the complexity of the languages they generate.

 $\hfill\square$ Type 0 \sim recursively enumerable.

There are no restrictions for the <u>rules in P</u>.

□ Type 1 ~ context-sensitive. For all rules $w \to w' \in P$ holds: $|w| \le |w'|$

$\hfill\square$ Type 2 \sim context-free.

For all rules $w \to w' \in P$ holds: w is a single variable; i.e., $w \in N$.

$\hfill\square$ Type 3 \sim regular.

The grammar is of Type 2, and for all rules, $w \to w'$ additionally holds: $w' \in (\Sigma \cup \Sigma N)$, i.e., the right-hand sides of the rules consist either of a terminal symbol or of a terminal symbol followed by a non-terminal.



Definition 3 (Language of Type)

A language $L \subseteq \Sigma^*$ is called a Type 0 (Type 1, Type 2, Type 3) language if there exists a Type 0 (Type 1, Type 2, Type 3) grammar *G* such that L(G) = L.

- $\hfill\square$ The Chomsky hierarchy represents a hierarchy with genuine subset relationships: Type 3 \subset Type 2 \subset Type 1 \subset Type 0
- □ All languages of types 1, 2, or 3 are decidable:
 - \rightarrow The <u>decision problem</u> if a word belongs to a language is decidable for all languages of Type 1, 2, or 3.
 - \rightarrow There is an algorithm that, given a grammar G and a word w, decides in finite time if $w \in L(G)$ or not.
- The set of Type 0 languages is identical to the set of recursively enumerable or semi-decidable languages. Therefore, there are Type 0 languages that are not decidable.

There are countably infinite grammars for generating a recursively enumerable language (Type 0)

□ In compiler theory and natural language processing, Type 3 languages and grammars (lexical analysis, tokenization) and Type 2 languages and grammars (syntactic structure analysis) play a central role.

Calculi for Regular Languages

Different calculi for generating words in a regular language:

- (a) finite acceptor or automaton
- (b) regular expressions

(d) Specifying finitely many equivalence classes (of the <u>Nerode relation</u>)

(c) Type 3 grammar

Calculi for Regular Languages

Different calculi for generating words in a regular language:

- (a) finite acceptor or automaton
- (b) regular expressions

(d) Specifying finitely many equivalence classes (of the <u>Nerode relation</u>)

(c) Type 3 grammar



Summary

	Language generation calculi
Туре 0	Type 0 grammar
	Turing machine
Туре 1	context-senstive grammar
	linear bounded Turing machine [Wikipedia]
Type 2	context-free grammar
	pushdown automaton
Туре З	regular grammar (Type 3 grammar)
	deterministic/non-deterministic finite automaton
	regular expression

Summary

	Language generation calculi		Complexity of the word problem [Wikipedia]
Туре 0	Type 0 grammar	Туре 0	undecidable
	Turing machine	Type 1	exponential complexity, NP-hard
Type 1	context-senstive grammar	Type 2	$O(n^3)$
	linear bounded Turing machine [Wikipedia]	Туре З	linear complexity
Type 2	context-free grammar pushdown automaton		
Туре 3	regular grammar (Type 3 grammar) deterministic/non-deterministic finite automaton regular expression		

Chapter NLP:I

- I. Natural Language Processing Basics
 - □ Grammars
 - Regular Expressions

Regular Expressions [Kastens] Syntax

A regular expression R can be composed recursively as follows. F and G denote regular expressions.

	R	Semantics
1.	a	The letter a
2.	FG	Concatenation
3.	$F \mid G$	Alternation
4.	(F)	Grouping
5.	F^+	Non-empty sequence of words from $L(F)$
6.	F^*	Arbitrary sequence of words from $L(F)$
7.	F^n	Sequence of n words from $L(F)$
8.	arepsilon	The empty word

Regular Expressions [Kastens] Syntax

A regular expression R can be composed recursively as follows. F and G denote regular expressions.

	R	Language $L(\mathbf{R})$	Semantics
1.	a	$\{a\}$	The letter a
2.	FG	$\{fg \mid f \in L(F), g \in L(G)\}$	Concatenation
3.	$F \mid G$	$\{f \mid f \in L(F)\} \cup \{g \mid g \in L(G)\}$	Alternation
4.	(F)	(L(F))	Grouping
5.	F^+	$\{f_1 f_2 \dots f_n \mid f_i \in L(F), n \ge 1, i = 1, \dots, n\}$	Non-empty sequence of words from $L(F)$
6.	F^*	$\{\varepsilon\} \cup L(F^+)$	Arbitrary sequence of words from $L(F)$
7.	F^n	$\{f_1 f_2 \dots f_n \mid f_i \in L(F), \ i = 1, \dots, n\}$	Sequence of n words from $L(F)$
8.	ε	$\{\varepsilon\}$	The empty word



Every word in the language of this regular expression consists of two ones, followed by any number of zeros or ones, followed by two zeros.

Regular Expressions [Kastens]

Examples (continued)

R	Name of language $L(\mathbf{R})$	Words from $L(\mathbf{R})$
(a b) (c d ε)	Abc	ac, bc, ad, bd, a, b
$Dear(\varepsilon \mid est)$ (Sir Madam)	Salutation	Dear Sir

Regular Expressions [Kastens]

Examples (continued)

R	Name of language $L(\mathbf{R})$	Words from $L(\mathbf{R})$
(a b) (c d ε)	Abc	ac, bc, ad, bd, a, b
$Dear(\varepsilon \mid est)$ (Sir Madam)	Salutation	Dear Sir
0 1 9	Digit	7
a b z	sLetter	х
A B Z	cLetter	В
sLetter cLetter	Letter	m, N

Regular Expressions [Kastens]

Examples (continued)

R	Name of language $L(\mathbf{R})$	Words from $L(\mathbf{R})$
(a b) (c d ε)	Abc	ac, bc, ad, bd, a, b
$Dear(\varepsilon \mid est)$ (Sir Madam)	Salutation	Dear Sir
0 1 9	Digit	7
a b z	sLetter	X
A B Z	cLetter	В
sLetter cLetter	Letter	m, N
Letter (Letter Digit)*	Label	Maximum, min7, a
Digit ⁺ .Digit ²	MoneyAmount	23.95, 0.50
(cLetter cLetter ² cLetter ³)– (cLetter cLetter ²)– (Digit Digit ² Digit ³ Digit ⁴)	LicensePlatesDE	PB-AS-0815
1 ³ (1 0)* 0 ³	Dual	1111000, 1111101010000

Examples (continued)

An important use of regular expressions in languages used for natural language processing is the specification of text patterns.

Example: Display of all file names of the form

```
"winter-term(0|1|2|3|4|5|6|7|8|9)<sup>2</sup>.html"
```

□ Unix-Shell.

```
ls winter-term[0-9][0-9].html
```

```
D PHP.
$d = "[0-9]";
preg match("/winter-term$d$d\.html/", $files)
```

 If names of regular expressions are used in other regular expressions, they must be identified as part of the meta language.

Here: Use of italics.

- □ Each scripting language for text processing uses a different syntax to specify regular expressions; the construction principles and power are comparable.
- □ The specification of regular expressions in PHP is taken from the Perl scripting language.

Definition 4 (Exact Matching Problem [Gusfield 1997])

Given a string R called the pattern and a longer string T called the text, the exact matching problem is to find all occurrences, if any, of pattern R in text T.

Example:

- \Box Given R = aba and T = bbabaxababay, then R occurs in T starting at locations 3, 7, and 9.
- \Box Two occurrences of *R* overlap.

Importance of the exact matching problem:

- □ Exact matching is solved for searching local files and data.
- □ Main bottleneck is the speed with which the data can be read into memory.
- \Box Speed rests on "preprocessing" pattern R (or text T) to enable skipping during search.

Perl Compatible Regular Expressions (PCRE) [pcre.org, PHP, Python, JavaScript, Java, .NET]

Many programming languages use a Perl-like syntax for regular expressions:

"Delimiter Regular_Expression Delimiter [Flags]"

Pattern

- □ The delimiter must be a non-alphanumeric character.
- Optional flags influence the matching strategy. [pcre.org, PHP, Python, JavaScript, Java, .NET]
- Search and replace functions with regular expressions. [pcre.org, PHP, Python, JavaScript, Java, .NET]

Perl Compatible Regular Expressions (PCRE) [pcre.org, PHP, Python, JavaScript, Java, .NET]

Many programming languages use a Perl-like syntax for regular expressions:

"Delimiter Regular_Expression Delimiter [Flags]"

Pattern

- □ The delimiter must be a non-alphanumeric character.
- Optional flags influence the matching strategy. [pcre.org, PHP, Python, JavaScript, Java, .NET]
- Search and replace functions with regular expressions. [pcre.org, PHP, Python, JavaScript, Java, .NET]



Perl Compatible Regular Expressions (continued) [Regular Expressions (Syntax)]

A regular expression R can be composed recursively as follows. Let F and G denote regular expressions and L(F), L(G) their languages over the alphabet Σ .

	R [pcre.org]	Semantics
1.	$\begin{array}{c} a \\ \vdots \\ [\Sigma'] \\ [^{\Sigma'}] \\ \Sigma' \end{array}$	The letter <i>a</i> Any symbol from Σ (except line break) Character class: Any symbol from $\Sigma' \subseteq \Sigma$ No symbol from character class $\Sigma' \subseteq \Sigma$ Escaped letter from a controlled $\Sigma' \subset \Sigma$ denoting a character class
2. 3. 4.	$FG \\ F \mid G \\ (F)$	Concatenation Alternation Grouping or capturing
5. 6. 7.	$F?\ F+\ F*\ F\{n\}\ F\{m,n\}$	<i>F</i> is optional (same semantics as $F \mid \varepsilon$) Non-empty sequence of words from $L(F)$ Arbitrary sequence of words from $L(F)$ Sequence of <i>n</i> words from $L(F)$ Sequence with at least <i>m</i> and at most <i>n</i> words from $L(F)$
	^ or \$	Assertion: Beginning or end of the text string (nothing may precede or follow it)

Regular Expressions PCRE: Character Classes [pcre.org]

Character classes match a single token from a list of characters.

R	Semantics	Example	
[]	Positive character class matches listed character	[aeiou] matches vowels	
[^]	Negative character class matches only unlisted characters	[^aeiou] matches anything except for vowels	
[x - y]	Range: matches any letter between x and y inclusive	[a-zA-Z] matches alphabetic characters	
\d	Matches any decimal digits. Short for [0-9]		
\D	Matches all except decimal digits. Short for [^0-9]		
\s	Matches white space characters. Short for $[\t\n\r\f\v]$		
\S	Matches all except white spaces. Short for $[^tt n\r]$		
\w	Matches alphanumeric characters. Short for [a-zA-Z0-9]		
$\setminus W$	Matches non-alphanumeric characters. Short for [^a-zA-Z0-9]]	

- □] needs to be escaped inside character classes to be matched: [^]] does not match a but, e.g.,]].
- □ The range in a character class denotes a series of <u>codepoints</u>. As such, [--2] is valid and identical to [-./012]. [A-z] matches not only alphabetic characters.
- **D** Character classes cannot be nested but \d, \D, \s, \dots can be used inside character classes:
 - $[^[a-z][A-Z]]$ matches the regex $[^[a-z], [A-Z]]$, and] in sequence.
 - $[^{\s}w]$ matches any character that is neither a space, nor a word-character.

PCRE: Groups and Backreferencing [pcre.org]

Groups give precedence to subformulas and can store (capture) matches for later reference.

R	Semantics	Example
(<i>F</i>)	Capture Group: Groups F and stores the matched value	(ab) + matches ab, abab,
(?:F)	Non-capturing group: Groups F	(?:ab) + matches ab, abab,
(?< <i>name</i> > <i>F</i>)	Named capture group: Groups F and stores it under a na	me
$\setminus n$	Backreference the <i>n</i> -th capture group	$(\w+)\1$ matches fully reduplicated words
\k{ name }	Backreference capture group named "name"	$(? < word > \w+) \k \{word\}$ is the same as above

Capture groups store the last match.

The regex ([ab]) + 1 matches abb but not aba.

□ Named capture groups can extract information from different positions in the regex.

The regex $d\d.(?<month>dd).dd(?<month>dd)/dddextracts the month from 16.01.70 and 01/16/70.$

- \Box The n syntax is highly ambiguous:
 - If *n* is an octal number ($\backslash 1, \backslash 2, \ldots, \backslash 7, \backslash 10, \ldots$), it matches the character point with value *n*.
 - If n is a single non-octal digit ($\setminus 8, \setminus 9$), it matches the respective character.
 - If *n* is a valid index for a capture group, the above is irrelevant and the capture group is backreferenced.

 $\8$ matches the character 8 but () () () () () () () $\8$ only matches the empty word. $\49$ gives an error because there is no 49th capture group but $\50$ matches (.

- → PCRE introduces gn or $g\{n\}$ to reference capture groups and n for octal codepoints. These are less common in other flavors of regex.
- □ Regular expressions compatible with PCRE2 are more powerful than regular grammars. The language $L = \{ww \mid w \in \{0,1\}^*\}$ is not context free but can be decided by the regex ([01]) *\1.

Quantifiers allow for repeated matching of a regular expression F.

R	Semantics	Example
F { m }	Sequence of m words from $L(F)$	\d{2} matches exactly two digits
$F\{m,\}$	Sequence of at least m words from $L(F)$	\d{2, } matches two or more digits
$F\{m,n\}$	Sequence of at least m and at most n words from $L(F)$	\d{2,5} matches two to five digits
F { , n }	Sequence of at most n words from $L(F)$ (since PCRE2 10.43)	\d{,10} matches up to ten digits
F?	Equivalent to $F\{0,1\}$	\w? matches at most one word-character
$F\star$	Equivalent to $F\{0, \}$	\w* matches any number of word-characters
F+	Equivalent to $F\{1,\}$	$\web w$ + matches one or more word-characters

Default. Greedy matching: as much as possible and backtracking until a match is found. $d*\d$ tries to match 1234

Quantifiers allow for repeated matching of a regular expression F.

R	Semantics	Example
$F \{m\}$	Sequence of m words from $L(F)$	\d{2} matches exactly two digits
$F\{m,\}$	Sequence of at least m words from $L(F)$	\d{2, } matches two or more digits
$F\left\{m,n ight\}$	Sequence of at least m and at most n words from $L(F)$	\d{2,5} matches two to five digits
F { , n }	Sequence of at most n words from $L(F)$ (since PCRE2 10.43)	\d{,10} matches up to ten digits
F? F* F+	Equivalent to $F\{0, 1\}$ Equivalent to $F\{0, \}$ Equivalent to $F\{1, \}$	<pre>\w? matches at most one word-character \w* matches any number of word-characters \w+ matches one or more word-characters</pre>

Default. Greedy matching: as much as possible and backtracking until a match is found.
 \d*\d\d tries to match 1234 but \d\d cannot be matched

Quantifiers allow for repeated matching of a regular expression F.

R	Semantics	Example
$F \{m\}$	Sequence of m words from $L(F)$	\d{2} matches exactly two digits
$F\{m,\}$	Sequence of at least m words from $L(F)$	\d{2, } matches two or more digits
$F\left\{m,n ight\}$	Sequence of at least m and at most n words from $L(F)$	\d{2,5} matches two to five digits
F { , n }	Sequence of at most n words from $L(F)$ (since PCRE2 10.43)	\d{,10} matches up to ten digits
F?	Equivalent to $F\{0,1\}$	\w? matches at most one word-character
$F\star$	Equivalent to $F\{0, \}$	\w* matches any number of word-characters
F+	Equivalent to $F\{1, \}$	$\web w$ + matches one or more word-characters

Default. Greedy matching: as much as possible and backtracking until a match is found.
 \d*\d\d tries to match 1234 but \d cannot be matched

Quantifiers allow for repeated matching of a regular expression F.

R	Semantics	Example
$F\left\{m ight\}$	Sequence of m words from $L(F)$	\d{2} matches exactly two digits
$F\{m,\}$	Sequence of at least m words from $L(F)$	\d{2, } matches two or more digits
$F\left\{m,n ight\}$	Sequence of at least m and at most n words from $L(F)$	\d{2,5} matches two to five digits
F { , n }	Sequence of at most n words from $L(F)$ (since PCRE2 10.43)	\d{,10} matches up to ten digits
F? F* F+	Equivalent to $F\{0, 1\}$ Equivalent to $F\{0, \}$ Equivalent to $F\{1, \}$	<pre>\w? matches at most one word-character \w* matches any number of word-characters \w+ matches one or more word-characters</pre>

□ Default. Greedy matching: as much as possible and backtracking until a match is found. \d*\d\d tries to match 1234 and but \d\d matches as well \rightarrow success!

Quantifiers allow for repeated matching of a regular expression F.

R	Semantics	Example
$F \{ m \}$	Sequence of m words from $L(F)$	\d{2} matches exactly two digits
$F\{m_{r}\}$	Sequence of at least m words from $L(F)$	\d{2, } matches two or more digits
$F\{m,n\}$	Sequence of at least m and at most n words from $L(F)$	\d{2,5} matches two to five digits
F { , n }	Sequence of at most n words from $L(F)$ (since PCRE2 10.43)	\d{,10} matches up to ten digits
F? F*	Equivalent to $F\{0,1\}$ Equivalent to $F\{0,\}$	<pre>\w? matches at most one word-character \w* matches any number of word-characters</pre>
F+	Equivalent to $F\{1, \}$	$\web w$ + matches one or more word-characters

□ Modify the matching behavior by adding ? (lazy) or + (possessive) behind the quantifier.

\d+? matches as few digits as possible, but at least one. \d++ matches greedily without backtracking

Strings		R		
	".*"	".*?"	".*+"	" [^ "] * + "
"hi"	{"hi"	} {"hi"}	Ø	{"hi"}
"foo"	"bar" {"foo	" "bar"} {"foo","bar"}	Ø	{"foo","bar"}

- D Possessive matching is more efficient than greedy or lazy matching because no backtracking is performed.
- □ Possessive matching is a special case of an atomic group (written (?>F)). Atomic groups cannot be backtracked.

The regex (?>a|aa) a matches aa but not aaa because for the latter it would initially match a for the capture group (the first matching option in the alternative) and then have to backtrack when the end of the string is reached.

What strings would (?>aa|a) a match?

 \Box F?+, F*+, and F++ are short for (?>F?)+, (?>F*)+, and (?>F+)+ respectively.

PCRE: Assertions [pcre.org]

Assertions test if a regex would match without matching it.

R	Semantics	
(?=F) (?!F) (?<=F) (? F)</td <td colspan="2"> (Positive) lookahead: Asserts that F matches at the current location. Negative lookahead: Asserts that F does not match at the current location. (Positive) lookbehind: Asserts that F matches right before the current location. Negative lookbehind: Asserts that F does not match right before the current location. (?<=\s) (?<!--\$)\d+ matches any number in the text that is not preceded by a \$.</li--> </td>	 (Positive) lookahead: Asserts that F matches at the current location. Negative lookahead: Asserts that F does not match at the current location. (Positive) lookbehind: Asserts that F matches right before the current location. Negative lookbehind: Asserts that F does not match right before the current location. (?<=\s) (?<!--\$)\d+ matches any number in the text that is not preceded by a \$.</li--> 	
^ \$	Begin of string anchor: Asserts the start of the input sequence. End of string anchor: Asserts the end of the input sequence. ^[01]+\\$ checks that the input entirely represents a binary number	
\b \B	Asserts a word boundary. Short for $(? <= \setminus W)$ $(? = \setminus W)$ $(? <= \setminus W)$ $(? = \setminus W)$ Asserts the opposite of $\setminus b$. Short for $(? <= \setminus W)$ $(? = \setminus W)$ $(? <= \setminus W)$ $(? = \setminus W)$ $\setminus b[Tt]he \setminus b$ matches the word "the" but not if it occurs inside another word as in "atheist". $\setminus Bing \setminus b$ matches the suffix -ing. For example in "climbing" or "sing" but not in "singer".	

- \Box ^, \$, \b, \B are simple assertions and (?=*F*), (?!*F*), (?<=*F*), (?<!*F*) are assertion groups.
- □ Simple assertions cannot be followed by quantifiers.

For example, b+ is illegal but (b) + is allowed.

NLP:I-44 Natural Language Processing Basics