



Java™ String and Text Processing

Lecturer: Nataliia Stulova

Teaching assistant: Mohammadreza Hazirprasand

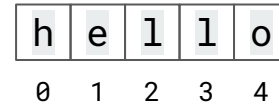
Software Composition Group
University of Bern
14 October 2020

Strings



java.lang.String

The `String` class represents character strings.



- Java strings are objects
- Java strings are **immutable** - once created, cannot be changed
- BUT you can query and search for substrings

```
String s = "hello";  
char c = s.charAt(1);           // c == 'e'  
int i = s.indexOf(l);          // i == '2'  
String sub = s.substring(1,3); // sub == "el"
```



Special characters

<code>\'</code>	Single quotes. Only need to escape if a string is wrapped in Single quote as well.
<code>\t</code>	Tab
<code>\b</code>	Backspace
<code>\r</code>	Return
<code>\f</code>	Formfeed
<code>\n</code>	Newline



Pretty printing

- control how non-string values appear in a string: decimal point digits number, time and date, indentation
- useful when you work with templated text messages

```
String fs;  
fs = String.format("The value of the float variable is " +  
                  "%f, while the value of the integer " +  
                  "variable is %d, and the string " +  
                  "is %s", floatVar, intVar, stringVar);  
System.out.println(fs);
```

Regular Expressions



Regular expressions

A regular expression (shortened as regex, regexp) is a sequence of characters that define a *search pattern*. Usually such patterns are used by string-searching algorithms for "find" or "find and replace" operations on strings, or for input validation.

Use: in search engines, search and replace dialogs of word processors and text editors, in text processing utilities such as `sed` and `awk` and in lexical analysis.



java.util.regex API

There are many different flavors to choose from, such as `grep`, `Perl`, `Tcl`, `Python`, `PHP`, and `awk`. The regular expression syntax in the `java.util.regex` API is most similar to that found in `Perl`.

The `java.util.regex` package primarily consists of three classes:

- `Pattern`,
- `Matcher`, and
- `PatternSyntaxException`.

This lecture covers the basics of <https://docs.oracle.com/javase/tutorial/essential/regex/intro.html>



regex: character classes

<code>[abc]</code>	a, b, or c (simple class): <code>[bc]at</code> matches <code>bat</code> and <code>cat</code>
<code>[^abc]</code>	Any character except a, b, or c (negation): <code>[^r]at</code> matches <code>fat</code> but not <code>rat</code>
<code>[a-zA-Z]</code>	a through z, or A through Z, inclusive (range)
<code>[a-d[m-p]]</code>	a through d, or m through p: <code>[a-dm-p]</code> (union)
<code>[a-z&&[def]]</code>	d, e, or f (intersection)
<code>[a-z&&[^bc]]</code>	a through z, except for b and c: <code>[ad-z]</code> (subtraction)
<code>[a-z&&[^m-p]]</code>	a through z, and not m through p: <code>[a-lq-z]</code> (subtraction)



regex: predefined character classes

<code>.</code>	Any character (may or may not match line terminators)
<code>\d</code>	A digit: <code>[0-9]</code>
<code>\D</code>	A non-digit: <code>[^0-9]</code>
<code>\s</code>	A whitespace character: <code>[\t\n\x0B\f\r]</code>
<code>\S</code>	A non-whitespace character: <code>[^\s]</code>
<code>\w</code>	A word character: <code>[a-zA-Z_0-9]</code>
<code>\W</code>	A non-word character: <code>[^\w]</code>



regex: quantifiers

<code>X?</code>	X, once or not at all
<code>X*</code>	X, zero or more times
<code>X+</code>	X, one or more times
<code>X{n}</code>	X, exactly n times
<code>X{n, }</code>	X, at least n times
<code>X{n, m}</code>	X, at least n but not more than m times



regex: boundary matches

<code>^</code>	The beginning of a line
<code>\$</code>	The end of a line
<code>\b</code>	A word boundary
<code>\B</code>	A non-word boundary
<code>\A</code>	The beginning of the input
<code>\G</code>	The end of the previous match
<code>\Z</code>	The end of the input but for the final terminator, if any
<code>\z</code>	The end of the input



java.util.regex API use

```
String regex = "a{2}+";  
Pattern pattern = Pattern.compile(regex);
```

regular expression string

```
String text = "aaaabbbbba";  
Matcher matcher = pattern.matcher(text);
```

text to search the pattern in

```
while (matcher.find()) {  
    System.out.println(  
        String.format(  
            "substring found at index [%d,%d]: %s",  
            matcher.start(),  
            matcher.end(),  
            matcher.group()));  
}
```

java.util.regex API use

```
String regex = "a{2}+";  
Pattern pattern = Pattern.compile(regex);
```

regular expression string

```
String text = "aaaabbbbba";  
Matcher matcher = pattern.matcher(text);
```

text to search the pattern in

```
while (matcher.find()) {  
    System.out.println(  
        String.format(  
            "substring found at index [%d,%d]: %s",  
            matcher.start(),  
            matcher.end(),  
            matcher.group()));  
}
```

```
substring found at index [0,2]: aa  
substring found at index [2,4]: aa
```

Practice



Exercise: regular expression patterns

write a `WordFilter.java` class with 5 static methods:

- `wordsAtLineBegin()` - find all words at line beginnings
- `wordsOfLength(int len)` - find all words of length `len` (e.g., 8 characters)
- `wordsAllCaps()` - find all words that consist of capital letters
- `wordsFirstCapital()` - find all words that start with a capital letter
- `wordsInBrackets()` - find all sub-sentences that are enclosed in `()` brackets

I/O

Use the example string (next slide) as input. Print results of each search to `System.out`.



Exercise: sample text

```
String inputText = "What separates API specifications from a programming guide are examples,\n"
    + "definitions of common programming terms, certain conceptual overviews (such as\n"
    + "metaphors), and descriptions of implementation bugs and workarounds. There is no\n"
    + "dispute that these contribute to a developer's understanding and help a\n"
    + "developer write reliable applications more quickly. However, because these do\n"
    + "not contain API \"assertions\", they are not necessary in an API specification.\n"
    + "You can include any or all of this information in documentation comments (and\n"
    + "can include custom tags, handled by a custom doclet, to facilitate it). At Java\n"
    + "Software, we consciously do not include this level of documentation in doc\n"
    + "comments, and instead include either links to this information (links to the\n"
    + "Java Tutorial and list of changes) or include this information in the same\n"
    + "documentation download bundle as the API spec -- the JDK documentation bundle\n"
    + "includes the API specs as well as demos, examples, and programming guides.";
```

- taken from <https://www.oracle.com/technical-resources/articles/java/javadoc-tool.html>
- use as the input text