

COMP284 Scripting Languages

Lecture 4: Perl (Part 3)

Handouts

# Assignment Project Exam Help

<https://eduassistpro.github.io>

Department of Computer

School of Electrical Engineering, Electronics, and

University of Liverpool

Add WeChat edu\_assist\_pro

# Assignment Project Exam Help

## ① Regul

Intr

Ch

Character classes

Quantifiers

<https://eduassistpro.github.io>

Add WeChat edu\_assist\_pr

# Regular expressions: Motivation

Suppose you are testing the performance of a new sorting algorithm by measuring its runtime on randomly generated arrays of numbers of a given length:

```
Generating an unsorted array with 10000 elements took 1.250 seconds
```

```
Sortin
```

```
Genera
```

```
Sortin
```

```
Genera
```

```
Sorting took 8.951 seconds
```

Your task is to write a program that determines the average runtime of the sorting algorithm:

```
Average runtime for 10000 elements is 8.886 seconds
```

Solution: The regular expression `/^Sorting took (\d+\.\d+) seconds/` allows us to get the required information

→ Regular expressions are useful for information extraction

# Regular expressions: Motivation

Suppose you have recently taken over responsibility for a company's website. You note that their HTML files contain a large number of URLs containing superfluous occurrences of `.`. . . e.g.

```
http://www.myorg.co.uk/info/refund/ ../vat.html
```

Your task is  
equivalent

```
http://
```

while making sure that relative URLs like

```
../videos/list.html
```

are preserved

Solution: `s![^\/]+/\.\.!!;` removes a superfluous dot-segment

↪ Substitution of regular expressions is useful for text manipulation

# Regular expressions: Introductory example

```
\Ahttps?:\\\/[^\\/]+\\.w\.\/(cat|dog)\\\/\1
```

## Assignment Project Exam Help

- \A is an **assertion or anchor**
- h, t, p, s, :, \/, c, a, t, d, o, g are **characters**
- ? and
- [^\\/
- . is a **me**
- (cat|dog) is **alternation** within a **capture group**
- \1 is a **backreference to a capture group**

<https://eduassistpro.github.io>

Add WeChat [edu\\_assist\\_pr](#)

# Pattern match operation

- To match a **regular expression** *regexpr* against the special variable `$_` simply use one of the expressions `/regexpr/` or `m/regexpr/`

This is called a **pattern match**

- `$_` is the **target string** of the pattern match

- In a **scal**

depen

```
if (/^Ahttps
... }
```

```
if (m/^Ahttps?:\/\/[^\./]+\.[^\./]+\/w\/(rat|log)\/\1\/
... }
```

# Regular expressions: Characters

The simplest **regular expression** just consists of a sequence of

- alphanumeric characters and
- non-alphanumeric characters escaped by a backslash

that matches exactly this sequence of characters occurring as a substring in the target

```
$_ = "ababcbcd"
if (/cbc/) { print "Match\n"} else { print "No match\n" }
```

Output:

Match

```
$_ = "ababcbcdcdde";
if (/dbd/) { print "Match\n"} else { print "No match\n" }
```

Output:

No match

## Regular expressions: Special variables

- Often we do not just want to know whether a regular expression matches a target string, but retrieve additional information

The special variable `$-` can be used to retrieve the start position of the match

Note that

- The special variable `$1` returns the first match after the match
- The special variable `$&` returns the match it

```
$_ = "abcdefghijklmnopqrstuvwxyz"
if (/cbc/) { print "Match found at position $-[0]: $&\n" }
```

Output:

```
Match found at position 4: cbc
```



# Regular expressions: Special escapes

There are various **special escapes** and **metacharacters** that match more than one character:

.	Matches any character except \n
\w	Matches a 'word' character (alphanumeric)
\W	Matches a non-word character
\s	Match a whitespace
\S	Match a non-whitespace
\d	Match a decimal digit
\D	Match a non-digit
\p{UnicodeProperty}	Match <i>UnicodeProperty</i> characters
\P{UnicodeProperty}	Match non- <i>UnicodeProperty</i> characters

# Regular expressions: Unicode properties

- Each **unicode character** has one or more **properties**, for example, which script it belongs to

`\p{UnicodeProperty}` matches all characters that have a particular property

- `\P{U}`

- Exam

Ara	
ASCII	ASCII characters
Currency_Symbol	Currency symbols
Digits	Digits in all scripts
Greek	Greek characters
Han	Chinese kanxi or Japanese kanji characters
Space	Whitespace characters

See <http://perldoc.perl.org/perluniprops.html> for a complete list

## Regular expressions: Character class

- A **character class**, a list of characters, special escapes, metacharacters and unicode properties enclosed in square brackets, matches any **single character** from within the class, for example, `[ad\t\n\-\09]`

- One m  
for exa

- A **caret** that is, it matches any **single character** that is **not** from within the class, for example, `[^01a-z]`

```
$_ = "ababcbcdcde";
if (/[bc][b-e][^bcd]/) {
    print "Match at positions $_-[0] to $_, $+[0]-1, ":_&\n";
}
```

Output:

```
Match at positions 8 to 10: cde
```

# Quantifiers

- The constructs for regular expressions that we have so far are not sufficient to match, for example, natural numbers of arbitrary size

Also, writing a regular expressions for, say, a nine digit number would be tedious

This is ma

<i>regex</i>	
<i>regexpr</i> <sup>+</sup>	Match <i>regexpr</i> 1 or more times
<i>regexpr</i> <sup>?</sup>	Match <i>regexpr</i> 1 or 0 tim
<i>regexpr</i> <sup>n</sup>	Match <i>regexpr</i> exactl
<i>regexpr</i> { <i>n</i> ,}	Match <i>regexpr</i> at least <i>n</i>
<i>regexpr</i> { <i>n</i> , <i>m</i> }	Match <i>regexpr</i> at least <i>n</i> but not more than <i>m</i> times

Quantifiers are greedy by default and match the longest leftmost sequence of characters possible

# Quantifiers

<i>regexpr</i> *	Match <i>regexpr</i> 0 or more times
<i>regexpr</i> +	Match <i>regexpr</i> 1 or more times
<i>regexpr</i> ?	Match <i>regexpr</i> 1 or 0 times
<i>regexpr</i> { <i>n</i> }	Match <i>regexpr</i> exactly <i>n</i> times
<i>regex</i>	
<i>regex</i>	

## Example

```
$_ = "Sorting took 10.486 seconds";
if (/\\d+\\.\\d+/) {
    print "Match at positions $_-[0] to $_+[0]-1\n";
$_ = "E00481370";
if (/[A-Z]0{2}(\\d+)/) {
    print "Match at positions $_-[1] to $_+[1]-1, ":"$_1\\n";
```

## Output:

```
Match at positions 13 to 18: 10.486
Match at positions 3 to 8: 481370
```

# Quantifiers

Example:

```
$_ = "E00481370";
```

```
if (/^(\d+)/) {  
    print "Match at positions $_[0] to $_[0]+1, " . $_&\n";  
}
```

Output:

```
Match at posit
```

- The reg
  - As the example illustrates, the regular expression
    - matches as early as possible
    - matches as many digits as possible
- ↪ quantifiers are greedy by default

# Revision

## Read

- Chapter 7: In the World of Regular Expressions
- Chapter 8: Matching with Regular Expressions

of

R. L. Schön

Learning Perl.

O'Reilly, 2011

- <http://perldoc.perl.org/perlre.html>
- <http://perldoc.perl.org/perlretut.html>
- <http://www.perlfect.com/articles/regextutor.shtml>