Basics	Wildcard and multipliers	Special characters	Negation	Other functions	Programming
0	000	000	00	000	000

Regular Expressions

SET09103 Advanced Web Technologies

School of Computing Napier University, Edinburgh, UK Module Leader: Uta Priss

Basics	Wildcard and multipliers	Special characters	Negation	Other functions	Programming
0	000	000	00	000	000

Outline

Basics

Wildcard and multipliers

Special characters

Negation

Other functions

Programming

Basics	Wildcard and multipliers	Special characters	Negation	Other functions	Programming
•	000	000	00	000	000

Character by character match



NO match NO match

Note: "i" at the end means "ignore case"

Basics	Wildcard and multipliers	Special characters	Negation	Other functions	Programming
0	●OO	000	00	000	000

Wildcard and multipliers

. stands for "any character".

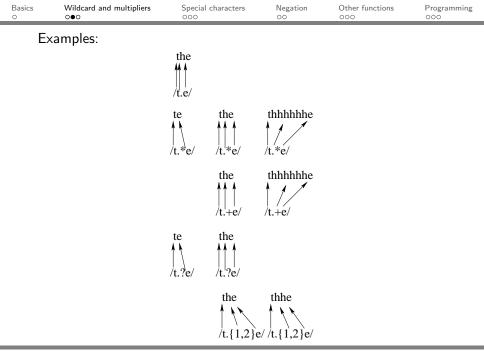
Basics	Wildcard and multipliers	Special characters	Negation	Other functions	Programming
0	•00	000	00	000	000

Wildcard and multipliers

. stands for "any character".

Multipliers:

+ stands for "at least one character"
* stands for "any number of characters (including 0)"
? stands for "at most one character" (i.e. either none or once)
{n,m} stands for "at least n times, at most m times"



Basics	Wildcard and multipliers	Special characters	Negation	Other functions	Programming
0	000	000	00	000	000

Exercise

What does /..\.19../ match: "12.1000" or "123.1900" or "12.2000"

Basics	Wildcard and multipliers	Special characters	Negation	Other functions	Programming
0	000	000	00	000	000

Exercise

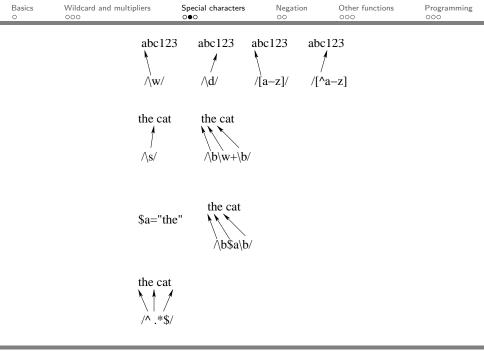
```
What does /..\.19../ match: 
"12.1000" or "123.1900" or "12.2000"
```

```
What does /hn*ell?o W...d/i match:
"Hello World" or "Hello Wood" or "Hell?o World"?
```

Basics	Wildcard and multipliers	Special characters	Negation	Other functions	Programming
0	000	•00	00	000	000

Special characters

$\setminus w$	word character (letter, digit or _)
[a-zA-Z]	letter
$\setminus W$	non-word character
[^a-zA-Z]	not a letter
d/	digit
$\setminus s$	space character (blank space, tab)
∖b	word boundary
^	beginning of line or string
\$	end of line or string



Basics	Wildcard and multipliers	Special characters	Negation	Other functions	Programming
0	000	000	00	000	000

Exercise

Which matches two consecutive words:

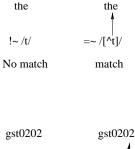
Basics	Wildcard and multipliers	Special characters	Negation	Other functions	Programming
0	000	000	O	000	000

Negation

 $word = \sim/[^a]/$ means that \$word must have one character which is not "a".

Basics	Wildcard and multipliers	Special characters	Negation	Other functions	Programming
0	000	000	00	000	000

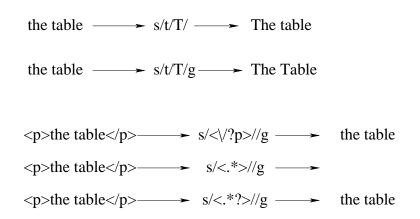
Examples



gst0202	gst0202;
	ţ
=~ /[^\w]/	$= \sim /[^w]/$
No match	match

Basics Wildcard and multipliers Speci	al characters Negation	Other functions P	rogramming
0 000 000	00	••••	00

Substitution



Basics	Wildcard and multipliers	Special characters	Negation	Other functions	Programming
0	000	000	00	000	000

Remembering patterns

Brackets are used for remembering patterns. The content of the first set of brackets can be retrieved with $\backslash 1.$ The second set of brackets with $\backslash 2,$ and so on.

Examples:

```
s/(the table)/\1/
/(.)\1/
<math>s/(.)(.)/(2)/(2)
```

Basics	Wildcard and multipliers	Special characters	Negation	Other functions	Programming
0	000	000	00	000	000

Split and Join (Implode)

```
$oldstring = "the,cat,sat,on,the,mat";
@array = split(/,/,$oldstring);
print @array;
# @array = ("the","cat","sat","on","the","mat")
$newstring = join(" ",@array);
# $newstring ="the cat sat on the mat"
```

Basics	Wildcard and multipliers	Special characters	Negation	Other functions	Programming
0	000	000	00	000	•00

Strategies

Instead of using one complicated regular expression, it is sometimes easier to use several simpler regular expressions combined with if statements.

For example: string starts with "a" and ends with "z":

if (
$$string = /^a.*z$$
)
if ($string = /^a/$ and $string = /z$)

Basics	Wildcard and multipliers	Special characters	Negation	Other functions	Programming
0	000	000	00	000	000

More Strategies

If a string needs to be processed ...

- ▶ from left to right, one character or one word at a time ⇒ split into array, then process array.
- ▶ from left to right, in some other regular manner ⇒ substr() can be used instead of regular expression.
- ► by checking whether some pattern exists ⇒ use regular expressions.

Basics	Wildcard and multipliers	Special characters	Negation	Other functions	Programming
0	000	000	00	000	000

Use of regular expressions in PHP

Searching:

if (preg_match("/the /i", \$line, \$matches)) {
 echo \$line,"
 matches: ",\$matches[0],"
";}

Replace:

```
$line = preg_replace("/T/", 't', $line);
```

Split:

```
$words = preg_split("/\s+/", $line);
```

Implode:

```
$newstring = implode(" ", $array);
```