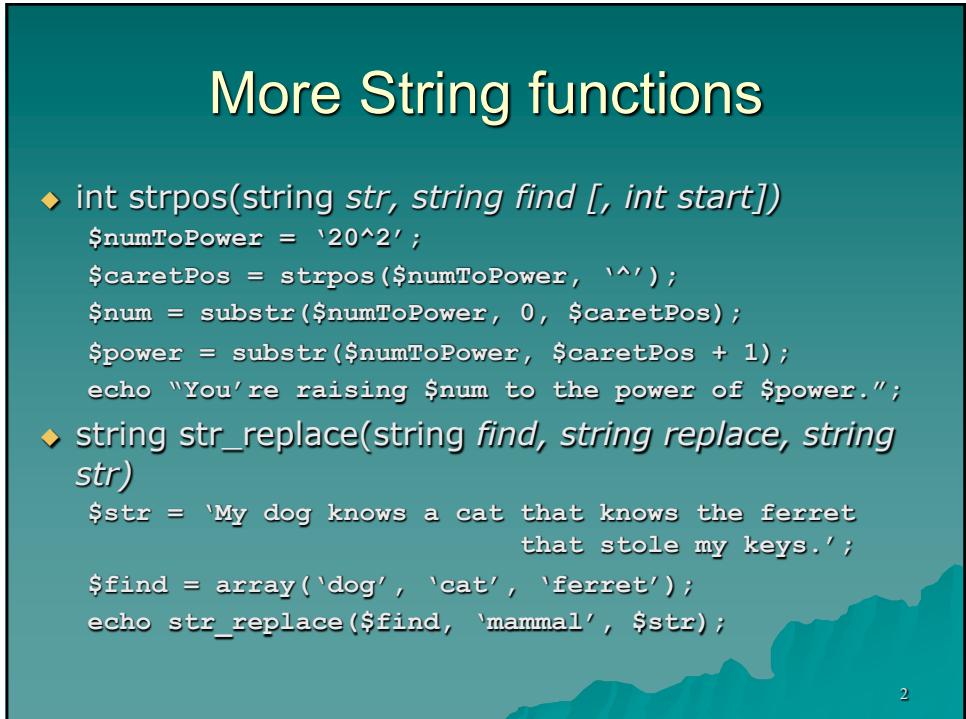




Web Development

Chapter 7. Regular Expressions

1



More String functions

- ◆ `int strpos(string str, string find [, int start])`

```
$numToPower = '20^2';
$caretPos = strpos($numToPower, '^');
$num = substr($numToPower, 0, $caretPos);
$power = substr($numToPower, $caretPos + 1);
echo "You're raising $num to the power of $power.";
```
- ◆ `string str_replace(string find, string replace, string str)`

```
$str = 'My dog knows a cat that knows the ferret
           that stole my keys.';
$find = array('dog', 'cat', 'ferret');
echo str_replace($find, 'mammal', $str);
```

2

2

1

Why regular expressions?

- ◆ Scripting problem may require:
 - verification of input from form
 - ◆ was input a 7 digit phone number
 - parsing input from a file
 - ◆ FirstName:LastName:Age:Salary
- ◆ PHP supports three pattern matching functions:
 - ereg(), split(), and ereg_replace()
- ◆ Regular expressions are used to define very specific match patterns

New version of PHP:

ereg → preg_match

split → preg_split

ereg_replace → preg_replace

3

3

ereg → preg_match

13.2. Switching From ereg to preg

13.2.1. Problem

You want to convert from using ereg functions to preg functions.

13.2.2. Solution

First, you have to add delimiters to your patterns:

```
preg_match('/pattern/', 'string')
```

For eregi() case-insensitive matching, use the /i modifier instead:

```
preg_match('/pattern/i', 'string');
```

When using integers instead of strings as patterns or replacement values,

```
$hex = dechex($number);
preg_match('/\x$hex/', 'string');
```

New version of PHP:

ereg → preg_match

split → preg_split

ereg_replace → preg_replace

4

4

2

The ereg() function

- ◆ Use ereg() to check if a string contains a match pattern:

```
$ret = ereg("search pattern", "target string");
```

↑
Set to 1 if pattern found. Set to 0 if not found.

↑
A set of normal or "special" characters to look for.

↑
A string value or string variable to search.

```
$ret = preg_match('/search pattern/', 'target string')
```

5

5

ereg() - example

- ◆ Consider the following

```
$name = 'Jake Jackson';
$pattern = 'ke';
if (ereg($pattern, $name)){
    print 'Match';
} else {
    print 'No match';
}
```

- ◆ This code outputs "Match" since the string "ke" is found.
- ◆ If \$pattern was "aa" the above code segment would output "No match"

6

6

3

Content

1. Regular Expression
2. Building an Example RE
3. Filter Input Data

7

7

Content

- ⇒ 1. Regular Expression
2. Building an Example RE
 3. Filter Input Data

8

8

1.1. What are regular expressions?

- ◆ Special pattern matching characters with specific pattern matching meanings.
 - Their meanings are defined by an industry standard (the IEEE POSIX 1003.2 standard).
 - For example, a caret symbol (^) returns a match when the pattern that follows starts the target string.

```
$part = 'AA100';
$pattern = '^AA';
if (ereg($pattern, $part)) {
    print 'Match';
} else {
    print 'No match';
}
```

Check if \$part starts with "AA"

Would be output if \$part was "AB100", "100AA", or "Apple".

9

9

1.2. Selected Pattern Matching Characters

Symbol	Description
^	<p><i>Matches when the following character starts the string.</i></p> <p>E.g. the following statement is <i>true</i> if \$name contains "Smith is OK", "Smithsonian", or "Smith, Black". It would be <i>false</i> if \$name contained only "SMITH" or "Smitty".</p> <pre>if (ereg('^Smith', \$name)) {</pre>
\$	<p><i>Matches when the preceding character ends the string.</i></p> <p>E.g. the statement below would be <i>true</i> if \$name contains "Joe Johnson", "Jackson", or "This is my son". It would be <i>false</i> if \$name contained only "My son Jake" or "MY SON".</p> <pre>if (ereg('son\$', \$name)) {</pre>

Slide 6a-10

10

1.2. Selected Pattern Matching Characters (2)

Symbol	Description
+	<i>Matches one or more occurrences of the preceding character.</i> For example, the statement below is <i>true</i> if \$name contains "AB101", "ABB101", or "ABBB101" is the right part". It would be <i>false</i> if \$name contained only "Part A101". <code>if (ereg('AB+101' , \$name)) {</code>
*	<i>Matches zero or more occurrences of the preceding character.</i> For example, the statement below is <i>true</i> if \$part starts with "A" and followed by zero or more "B" characters followed by "101", (for example, "AB101", "ABB101", "A101", or "A101 is broke"). It would be <i>false</i> if \$part contained only "A11". <code>if (ereg('^AB*101' , \$part)) {</code>
?	<i>Matches zero or one occurrences of the preceding character</i>

11

11

1.2. Selected Pattern Matching Characters (3)

Symbol	Description
.	<i>A wildcard symbol that matches any one character.</i> For example, the statement is <i>true</i> if \$name contains "Stop", "Soap", "Szxp", or "Soap is good". It would be <i>false</i> if \$name contained only "Sxp". <code>if (ereg('^S..p' , \$name)) {</code>
	<i>An alternation symbol that matches either character pattern.</i> For example, the statement below would be <i>true</i> if \$name contains "www.mysite.com", "www.school.edu", "education", or "company". It would be <i>false</i> \ if \$name contained only "www.site.net". <code>if (ereg('com edu' , \$name)) {</code>

Slide 6a-12

12

For example ...

- ◆ Regular expressions are case sensitive by default

```
<br>Enter product code (Use AB## format):<br>
    <input type="text" size="6" name="code">
<br> Please enter description:<br>
    <input type="text" size="50" name="description">
```

- ◆ Asks for a product code and description (not to contain "Boat" or "Plane").

13

13

A Full Script Example

- ◆ Consider an example script that enables end-user to select multiple items from a checklist.
 - A survey about menu preferences
 - Will look at how to send multiple items and how to receive them (later)

14

14

A Full Example ...

```
1. <html><head><title>Product Information Results </title>
2. </head><body>
3. <?php
4. $products = array('AB01'=>'25-Pound Sledgehammer',
                     'AB02'=>'Extra Strong Nails',
                     'AB03'=>'Super Adjustable Wrench',
                     'AB04'=>'3-Speed Electric Screwdriver');
5. if (ereg('boat|plane', $description)) {
6.     print 'Sorry, we do not sell boats or planes anymore';
7. } elseif (ereg('^AB', $code)) {
8.     if (isset($products["$code"])){
9.         print "Code $code Description: $products[$code]";
10.    } else {
11.        print 'Sorry, product code not found';
12.    }
13. } else {
14.     print 'Sorry, all our product codes start with "AB"';
15. } ?> </body></html>
```

Create a list of products.

Check if “boat” or “plane”.

Check if valid product number

15

15

The Output ...

The previous code can be executed at
<http://webwizard.aw.com/~phppqm/C6/drivSimple.html>



16

16

1.3. Using grouping characters

- Use **parentheses** to specify a group of characters in a regular expression.

Match Statement	Possible Matching Values
<pre>if (ereg('Dav(e id)', \$name)) {</pre>	"Dave", "David", "Dave was here"

- Above uses parentheses with “|” to indicate “Dav” can be followed by “e” or “id”.

Slide 6a-17

17

1.3. Using grouping characters (2)

- Now add in “^” and “\$” characters ...

Match Statement	Possible Matching Values
<pre>if (ereg('^ (d D)av(e id)\$', \$name)) {</pre>	"Dave", "David", "dave", "david"

Slide 6a-18

18

1.3. Using grouping characters (3)

- Use curly brackets to specify a range of characters
 - to look for a repeating of one or more characters
 - E.g.
 - L{3} matches 3 "L"s
 - L{3,} matches 3 or more "L"
 - L{2,4} matchs 2 to 4 "L"

Match Statements	Possible Matching Values
if (ereg('^L{3}\$', \$name)) {	"LLL" only
if (ereg('^L{3,}\$', \$name)) {	"LLL", "LLLL", "LLLLL", "LLLLL", and so on
if (ereg('^L{2,4}\$', \$name)) {	"LL", "LLL", or "LLLL" only

19

19

1.3. Using grouping characters (4)

- Use square brackets for character classes
 - to match one of character found inside them

Match Statement	Possible Matching Values
if (ereg('Sea[nt]!', \$name)) {	"Sean!", "Seal!", "Here comes Sean!"

Pattern: (r|f)at
Match string: carat

MATCH 1 - FINISHED IN 13 STEPS

```
1 //[(r|f)at/ |carat
2 //[(r|f)at/ |carat
3 //[(r|f)at/ |carat
4 //[(r|f)at/ |carat
5 //[(r|f)at/ |carat BACKTRACK
6 //[(r|f)at/ |carat
7 //[(r|f)at/ |carat
8 //[(r|f)at/ |carat BACKTRACK
9 //[(r|f)at/ |carat
10 //[(r|f)at/ |carat
11 //[(r|f)at/ |carat
12 //[(r|f)at/ |carat
13 //[(r|f)at/ |carat
# Match found in 13 steps.
```

Pattern: [rf]at
Match string: carat

MATCH 1 - FINISHED IN 7 STEPS

```
1 //|[rf]at/ |carat
2 //|[rf]at/ |carat
3 //|[rf]at/ |carat
4 //|[rf]at/ |carat
5 //|[rf]at/ |carat
6 //|[rf]at/ |carat
7 //|[rf]at/ |carat
# Match found in 7 steps.
```

Slide 6a-20

20

10

1.3. Using grouping characters (5)

- Use square brackets with range
 - More common to specify a range of matches
 - For example [0-9], [a-z] or [A-Z]

Match Statement	Possible Matching Values
<code>if (ereg('[0-9]', \$prodcode)) {</code>	“apple1”, “24234”, “suzy44”, “s1mple”

- Or use multiple characters at once ...

Match Statement	Possible Matching Values
<code>if (ereg('[A-Z][A-Z][0-9]', \$code)) {</code>	“AA9”, “Send product AZ9”, “MY12”

21

21

1.3. Using grouping characters (6)

- Using caret “^” and square brackets
 - When caret “^” is first character within square brackets it means “not”.

Match Statement	Possible Matching Values
<code>if (ereg('^5-9][0-9][A-Z]', \$code)) {</code>	“The AA9A is OK”, “Product 44X is down”, “It was 9Years ago.”

- Note: Within a character class, as in [^ . . .], “^” means *not*. Earlier saw how it can indicate that the character that follows the caret symbol *starts* the match pattern

22

22

1.4. Special Pre-defined character classes

Character Class	Meaning
<code>[:space:]</code>	<p><i>Matches a single space (Whitespace: newline, carriage return, tab, space, vertical tab) → [\n\r\t \x0B]</i></p> <p>E.g. the following matches if \$code contains “Apple Core”, “Alle y”, or “Here you go”; it does not match “Alone” or “Fun Time”:</p> <pre>if (ereg('e[:space:]', \$code)) {</pre>
<code>[:blank:]</code>	<p><i>Horizontal whitespace (space, tab) → [\t]</i></p>
<code>[:alpha:]</code>	<p><i>Matches any word character (uppercase or lowercase letters.).</i></p> <p>E.g., the following matches “Times”, “Treaty”, or “timetogo”; it does not match “#%^&”, “time” or “Time to go”:</p> <pre>if (ereg('e[:alpha:]', \$code)) {</pre>

23

23

1.4. Special Pre-defined character classes (2)

Character Class	Meaning
<code>[:upper:]</code>	<p><i>Matches any single upper case character and not lower case → [A-Z]</i></p> <p>E.g., the following matches “Home” or “There is our Home”, but not “home”, or “Our home“:</p> <pre>if (ereg('[[upper:]]ome', \$code)) {</pre>
<code>[:lower:]</code>	<p><i>Matches any single lower case character and not upper case → [a-z]</i></p> <p>E.g. the following matches “home” or “There is our home”, but not “Home”, or “Our Home“:</p> <pre>if (ereg('[[lower:]]ome', \$code)) {</pre>
<code>[:alpha:]</code>	<i>Matches any single alphabetic characters (letters) → [a-zA-Z]</i>
<code>[:alnum:]</code>	<i>Matches any single alphanumeric characters (letters) → [[0-9a-zA-Z]]</i>

24

24

1.4. Special Pre-defined character classes (3)

Character Class	Meaning
<code>[:digit:]</code>	<i>Matches any valid numerical digit (that is, any number 0-9)</i> $\rightarrow [0-9]$ E.g., the following matches “B12abc”, “The B1 product is late”, “I won bingo with a B9”, or “Product B00121”; it does not match “B 0”, “Product BX 111”, or “Be late 1”: <code>if (ereg('B[:digit:]', \$code)) {</code>
<code>[:punct:]</code>	<i>Matches any punctuation mark</i> $\rightarrow [-!"#$%&'()*+,./;=>?@[\\"\\]^_{}~]$ E.g., the following matches “AC101!”, “Product number.”, or “!!”, it does not match “1212” or “test”: <code>if (ereg('[:punct:]\$', \$code)) {</code>

25

25

1.4. Special Pre-defined character classes (4)

Character Class	Meaning
<code>[:<:]</code>	<i>Matches when the following word starts the string.</i>
<code>[:>:]</code>	<i>Matches when the preceding word ends the string</i>

E.g.,

```
// returns false  
ereg('[:<:]gun[:>]', 'the Burgundy exploded');  
// returns true  
ereg('gun', 'the Burgundy exploded');
```

26

26

Content

1. Regular Expression
2. Building an Example RE
3. Filter Input Data

27

27

2. Building an example RE

- ◆ Building Regular expressions is best done incrementally
- ◆ Lets look at a process to build a regular expression to validate a date input field:
 - mm/dd/yyyy format (for example, 01/05/2002 but not 1/5/02).

28

28

14

2.1. Determine the precise field rules

- ◆ What is valid input and invalid input
 - You might decide to allow 09/09/2002 but not 9/9/2002 or Sep/9/2002 as valid date formats.
- ◆ Work through several examples as follows:

Rule	Reject These
1. Only accept "/" as a separator	05 05 2002—Require slash delimiters
2. Use a four-digit year	05/05/02—Four-digit year required
3. Only date data	The date is 05/05/2002—Only date fields allowed 05/05/2002 is my date—Only date fields allowed
4. Require two digits for months and days	5/05/2002—Two-digit months required 05/5/2002—Two-digit days required 5/5/2002—Two-digit days and months required

29

29

2.2. Get the form and form-handling scripts working

- Build the input form and a “bare bones” receiving script
- For example: receives input of 1 or more characters:

```
if (ereg('.+', $date)) {
    print "Valid date= $date";
} else {
    print "Invalid date= $date";
}
```

Slide 6a-30

30

15

2.3. Start with the most specific term possible

- You know must have 2 slashes between 2 character month, 2 character day and 4 character year
- So change receiving script to:

```
if ( ereg( '.../.../....', $date ) ) {
    print "Valid date= $date";
} else {
    print "Invalid date= $date";
}
```

- So 12/21/1234 and fj/12/ffff are valid, but 1/1/11 is not.

31

31

2.4. Anchor the parts you can

- Add the “^” and “\$” quantifiers where possible.
- Also, can add the [[:digit:]] character class to require numbers instead of any character.
- So change receiving script to:

```
$two='[[:digit:]]{2}';
if ( ereg("^$two/$two/$two$two$", $date) )
{
    print "Valid date= $date";
} else {
    print "Invalid date= $date";
}
```

- So 01/16/2003, 09/09/2005, 01/12/1211, and 99/99/9999 are valid dates.

32

32

2.5. Get more specific if possible

- You might note that three more rules can be added:
 - The first digit of the month can be only 0, or 1. For example, 25/12/2002 is clearly illegal.
 - The first digit of a day can be only 0, 1, 2, or 3. For example, 05/55/2002 is clearly illegal.
 - Only allow years from this century allowed. Don't care about dates like 05/05/1928 or 05/05/3003.

```
$two='[:digit:]{2}';  
$month='[0-1][:digit:]';  
$day='[0-3][:digit:]';  
$year="2[:digit:]"$two";  
if ( ereg("^($month)/($day)/($year)$", $date ) ) {
```

Now input like
09/99/2001 and
05/05/4000 is illegal.

33

33

A Full Script Example

- ◆ Consider an example script that asks end-user for a date
 - Use regular expressions to validate
 - Use the following HTML input

```
<input type="text" size="10" maxlength="10" name="date">
```

34

34

A Full Example ...

```
1. <html>
2. <head><title>Decisions</title></head>
3. <body>
4. <?php
5.     $two='[[:digit:]]{2}';
6.     $month='[0-3][[:digit:]]';
7.     $day='[0-3][[:digit:]]';
8.     $year="2[[:digit:]]$two";
9.     if ( ereg("^($month)/($day)/($year)$", $date ) ) {
10.         print "Got valid date=$date <br>";
11.     } else {
12.         print "Invalid date=$date";
13.     }
14.?> </body></html>
```

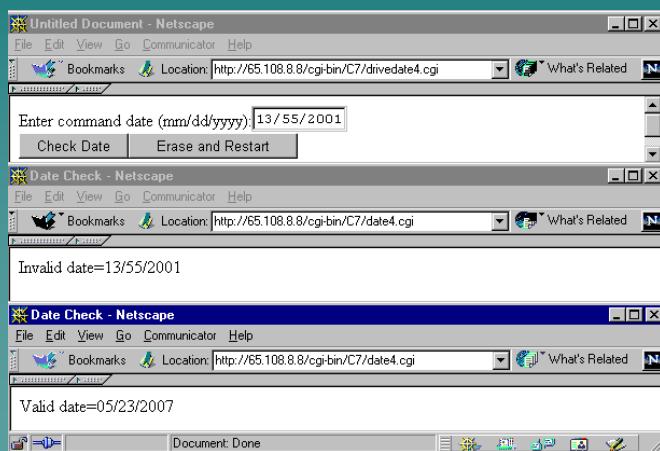
Use same regular expression as before

35

35

The Output ...

The previous code can be executed at
<http://webwizard.aw.com/~phppqm/C6/drivedate4.cgi>



36

36

18

Content

1. Regular Expression
2. Building an Example RE
3. Filter Input Data

37

37

3.1. Matching Patterns With `split()`

- Use `split()` to break a string into different pieces based on the presence of a match pattern.

```
Array variable that  
will contain resulting  
matches.    $output = split(search_pattern, target_string, max_match);  
String pattern with regular  
expression to match.  
A string to search  
Maximum number  
of matches  
to make (optional).
```

```
preg_split ( string $pattern , string $subject [, int $limit = -1 [, int $flags = 0 ]] )
```

38

38

19

3.1. Matching Patterns With split()

- Consider another example:

```
$line = 'Baseball, hot dogs, apple pie' ;
$item = split( ',', $line );
print ("0=$item[0] 1=$item[1] 2=$item[2]");
```

- These lines will have the following output:

```
0=Baseball 1= hot dogs 2= apple pie
```

39

39

3.1. Matching Patterns With split()

- When you know how many patterns you are interested can use list() along with split():

```
line = 'AA1234:Hammer:122:12';
list($partno, $part, $num, $cost)
    = split(':', $line, 4);
print "partno=$partno part=$part num=$num
cost=$cost";
```

- The above code would output the following:

```
partno=AA1234 part=Hammer num=122 cost=12
```

40

40

Example of split()

- ◆ As an example of split() consider the following:

```
$line = 'Please , pass      thepepper';
$result =  split( '[:space:]+', $line );
```

- ◆ Will results in the following:

```
$result[0] = 'Please' ;
$result[1] = ','
$result[2] = 'pass' ;
$result[3] = 'thepepper' ;
```

41

41

A Full Script Example

- ◆ Consider an example script that updates the date checker just studied:
 - Uses split() to further refine date validation
 - Uses the same input form:

```
<input type="text" size="10" maxlength="10" name="date">
```

42

42

A Full Example ...

```
1. <html>
2. <head><title>Date Check</title></head>
3. <body>
4. <?php
5.     $two='[[:digit:]]{2}';
6.     $month='[0-3][[:digit:]]';
7.     $day='[0-3][[:digit:]]';
8.     $year="2[[:digit:]]$two";
9.     if ( ereg("^(?{$month})/($day)/($year)", $date) ) {
10.         list($mon, $day, $year) = split( '/', $date );
11.         if ( $mon >= 1 && $mon <= 12 ) {
12.             if ( $day <= 31 ) {
13.                 print "Valid date mon=$mon day=$day year=$year";
14.             } else {
15.                 print "Illegal day specified Day=$day";
16.             }
17.         } else {
18.             print "Illegal month specified Mon=$mon";
19.         }
20.     } else {
21.         print ("Invalid date format= $date");
22.     }
23. ?></body></html>
```

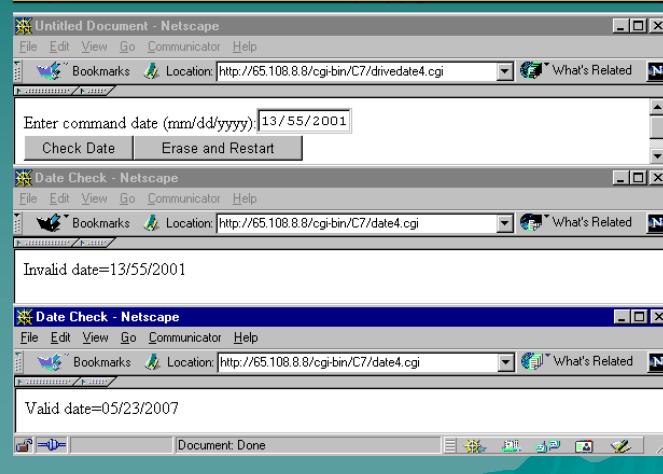
Use split() and list() to get month, day and year.

43

43

The Output ...

The previous code can be executed at
<http://webwizard.aw.com/~phppgm/C6/drivedate4.cgi>



44

44

3.2. Using ereg_replace()

- ◆ Use ereg_replace() when replacing characters in a string variable.
 - It can be used to replace one string pattern for another in a string variable.
 - E.g:

```
$start = 'AC1001:Hammer:15:150';
$end = ereg_replace('Hammer', 'Drill', $start );
print "end=$end";
```
 - The above script segment would output:
end=AC1001:Drill:15:150

45

45

Summary

- ◆ PHP supports a set of operators and functions that are useful for matching and manipulating patterns in strings:
 - The ereg() function looks for and match patterns
 - The split() function uses a pattern to split string values into as many pieces as there are matches.
 - The ereg_replace() function replaces characters in a string variable
- ◆ Regular expressions greatly enhance its pattern matching capabilities.

46

46