

Programování v PHP

Katedra softwarového inženýrství
Fakulta informačních technologií
České vysoké učení technické v Praze

© Pavel Štěpán, Helena Wallenfelsová, 2013

Regulární výrazy
BI-PHP



Regulární výrazy - testování, vyhledávání a úpravy řetězců dle masky (reg. výrazu) (Zde se pracuje s regulárními výrazy ve stylu jazyka Perl. Lze používat i styl POSIX)

Základní metaznaky regulárních výrazů:

<code>/^abc\$/</code>	<code>^</code> - začátek řetězce, <code>\$</code> - konec
<code>/^ab*c\$/</code>	<code>*</code> - opakování předchozího znaku nebo skupiny znaku 0,1,2,...
<code>/^ab+c\$/</code>	<code>+</code> - opakování 1,2,3,...
<code>/^ab?c\$/</code>	<code>?</code> - opakování 0,1 (nepovinné)
<code>/^ab{3,5}c\$/</code>	<code>{m,n}</code> - opakování m,m+1,m+2,...,n
<code>/^ab{3,}c\$/</code>	<code>{m,}</code> - opakování m, m+1, m+2, ...
<code>/^ab{3}c\$/</code>	<code>{m}</code> - opakování m-krát
<code>/^a(xyz)*c\$/</code>	<code>(...)</code> - skupina znaku pro opakování, podvýraz
<code>/^a(po ut st ct pa)c\$/</code>	<code>(...)</code> - kterákoliv z možností
<code>/^a[kxd]c\$/</code>	<code>[...]</code> - jeden z uvedených znaků
<code>/^a[^kxd]c\$/</code>	<code>[^...]</code> - jeden znak, různý od uvedených
<code>/^a[0-9]c\$/</code>	<code>[...-...]</code> - jeden znak ze zadaného intervalu
<code>/^a.c\$/</code>	<code>.</code> (tečka) - jeden libovolný znak

Některé zkratky:

<code>[0-9]</code>	<code>\d</code>
<code>[^0-9]</code>	<code>\D</code>
whitespace	<code>\s</code> (mezera, tabulátor, nový řádek, ...)
not whitespace	<code>\S</code>
<code>[A-Za-z0-9_]</code>	<code>\w</code>
<code>[^A-Za-z0-9_]</code>	<code>\W</code>
word boundary	<code>\b</code> (např. Java pro versus JavaScript)
not word boundary	<code>\B</code>

Příklady:

<code>/^[0-9]{3}[]?[0-9]{2}\$/</code>	PSC s nepovinnou mezerou
<code>/^\d{3}[]?\d{2}\$/</code>	totež (se zkratkami)
<code>/^[+-]?[1-9][0-9]*([,][0-9]+)?\$/</code>	číslo bez vedoucích nul s nepovinným značením a desetinnou částí
<code>/^ab+c\$/i</code>	<code>/.../img</code> (i - nerozlišovat malá/velká, m - multiline, g - global)
<code>/[\s]*,[\s]*/</code>	čárka, obklopená libovolným počtem mezer (i nulovým)
<code>/(['"])(['"])*\1/</code>	řetězec (neobsahující ' nebo ") obklopený znaky ' nebo " (\1 - řetězec, odpovídající 1. podvýrazu (zavorce))
<code>(\w+\.)*\w+@(\w+\.)+[A-Za-z]{2,3}</code>	mailová adresa (zjednodušená)
<code>Java(=?Script)</code>	pozitivní předbežný předpoklad (shoda jen, je-li nasledováno slovem Script; to se ale nepamatuje jako podvýraz)
<code>Java(?!Script)</code>	negativní předbežný předpoklad (Java není následováno Script, které se nepamatuje)

Funkce pro regulární výrazy, používané v tomto programu (detaily viz help):

`preg_match` - zjistuje shodu daného řetězce (2. param) s regulárním výrazem (1. param) a vrátí true/false. Může vrátit i pole pro shodu s různými podvýrazy

`preg_match_all` - nalezení všech shod (vrátí pole) - včetně podvýrazů

`preg_replace` - nahradí v řetězci všechny shody s reg.výrazem zadanou hodnotou

`preg_split` - rozdělí řetězec na podřetězce (vrátí pole) dle zadaného regulárního výrazu

```

<!DOCTYPE html>
<html>
    <head>
        <meta http-equiv="Content-Type" content="text/html; charset=UTF-8">
        <title>Jednoduchý program pro práci s regulárními výrazy</title>
    </head>
    <body>

        <?php
            $retez = $valRetez = "";
            $re = $valRe = "";
            $replSplit = $valReplSplit = "";
            $vystup = "";

            if (!empty($_GET["btnTest"]) || !empty($_GET["btnReplace"]) ||
                !empty($_GET["btnSplit"])) { // stisknuto tlačítko?
                $retez = $_GET["txtRetez"];
                $re = $_GET["txtVyzraz"];
                $replSplit = $_GET["txtReplSplit"];

                // Zapnuto vkládání backslashes? (novější verze PHP by neměly
                // používat magic quotes - automatické vkládání escape char)
                /*
                if (get_magic_quotes_gpc()) {
                    // stripslashes zruší přidání znaků \ (backslash)
                    $retez = stripslashes($retez);
                    $re = stripslashes($re);
                    $replSplit = stripslashes($replSplit);
                }
                */

                // str_replace použít k nahrazení případných úvozovek
                // v řetězci entitou &quot; - aby po vložení do tagu
                // input (jako value) nebyla úvozovka považována za
                // konec hodnoty u atributu value!! Např. value="/[;,]/"
                $valRetez = str_replace("'", '&quot;', $retez);
                $valRe = str_replace("'", '&quot;', $re);
                $valReplSplit = str_replace("'", '&quot;', $replSplit);
            }
        ?>

        <form action="re_test.php" method="get">
            Retezec:<br>
            <input type="text" name="txtRetez" size="80"
                value="<?php echo $valRetez;?>"><br>
            Reg. exp.:<br>
            <input type="text" name="txtVyzraz" size="80"
                value="<?php echo $valRe;?>"><br>
            Replace/Split:<br>
            <input type="text" name="txtReplSplit" size="80"
                value="<?php echo $valReplSplit;?>"><br>

            <input type="submit" name="btnTest" value="Test">
            <input type="submit" name="btnReplace" value="Replace">
            <input type="submit" name="btnSplit" value="Split"><br><br>
        </form>
    </body>
</html>

```

```

<?php
    if (!empty($_GET["btnTest"])){
        // testuje shodu
        // $vysledek = preg_match($re, $retez, $shody);

        // nalezne vsechny shody
        $vysledek = preg_match_all($re, $retez, $shody);

        if ($vysledek)
            echo "<b>Retezec odpovida regularnimu vyrazu</b><br>";
        else
            echo "<b>Retezec NEodpovida regularnimu vyrazu</b><br>";

        $vystup = "";
        foreach ($shody[0] as $shoda){ // pole vsech shod s celym re
            $vystup .= $shoda."<br>";
        }
    }
    elseif (!empty($_GET["btnReplace"])){
        // nahrazeni v retezci dle reg. vyrazu
        $vystup = preg_replace($re,$replSplit,$retez);

        // Priklad 1 - nahrazeni mezery, obklopené libovolným
        // počtem mezer carkou a mezerou
        // $text = "První , Druhý , Třetí";
        // echo "<br>".preg_replace("/[\s]*,[\s]*/"," ",$text)."<br>";

        // Priklad 2 - nahradit uvozovky nebo apostrofy, obklopující
        // řetězce (které neobsahující apostrof ani uvozovku)
        // hranatými závorkami
        // $text = "\"První\",'Druhý','Třetí'";
        // echo "$text<br>";
        // \1 - první nalezený podvzor = cast v závorkách
        // (pro regularní výraz)
        // $2 - druhý nalezený podvzor (pro nahrazování)
        // echo "<br>".preg_replace("/(['\"])([^\"]*)\\1/",
        // "[$2]",$text)."<br>";
    }
    elseif (!empty($_GET["btnSplit"])){
        // Příklad - rozdělení slov s libovolným počtem mezer
        // jako oddělovacem
        // $text = "První Druhý Třetí";
        // rozdeli řetězec podle regularního výrazu
        // $casti = preg_split("/[\s]+/",$text);
        // echo "<br>".implode(",",$casti)."<br><br>";

        // rozdělení řetězce podle regularního výrazu
        $casti = preg_split($replSplit,$retez);

        $vystup = "";
        foreach ($casti as $cast){
            $vystup .= $cast."<br>";
        }
    }
}
?>
<textarea name="taShody" rows="10" cols="80">
    <?php echo $vystup?></textarea>
</body>
</html>

```