

Dla większości prostych gramatyk można w łatwy sposób napisać wyrażenie regularne które będzie służyło do sprawdzania poprawności zdań z tą gramatyką.

Celem niniejszego laboratorium będzie zapoznanie się z wyrażeniami regularnymi (Regex) występującymi w języku C# oraz napisanie odpowiednich wyrażeń regularnych do sprawdzania poprawności zdań z zadanymi gramatykami.

## 1. Podstawy budowania wyrażeń regularnych (Regex)

Znaki:

. \$ ^ { [ ( | ) ] } \* + ? \

są w wyrażeniach regularnych znakami specjalnymi.

Przy użyciu poniższych symboli można zaznaczyć dopasowanie do 'klasy' znaków.

.	Dowolny znak
[abc]	Dowolny znak ze wskazanego zbioru
[^abc]	Dowolny znak różny od tych w zbiorze
[a-z0-9]	Dowolny znak z zakresu
\w	Dowolny znak typu: litera, cyfra, podkreślnik
\W	Dowolny znak inny niż: litera, cyfra, podkreślnik
\s	Dowolny biały znak (spacja, ...)
\S	Dowolny znak różny niż znak biały
\d	Dowolna cyfra
\D	Dowolny znak inny niż cyfra

Przy użyciu poniższych symboli można zaznaczyć miejsce w tekście gdzie ma dojść do dopasowania.

<code>^</code>	Dopasowanie musi się pojawić na początku tekstu
<code>\$</code>	Dopasowanie musi się pojawić na końcu tekstu
<code>\b</code>	Dopasowanie musi się znajdować na początku lub końcu słowa
<code>\B</code>	Dopasowanie nie może się znajdować na początku ani końcu słowa

Za pomocą poniższych symboli można określić ile razy dany znak (lub grupa znaków) ma wystąpić.

<code>*</code>	Zero lub więcej razy
<code>+</code>	Jeden lub więcej razy
<code>?</code>	Zero lub jeden raz
<code>{n}</code>	Dokładnie n razy
<code>{n,}</code>	Co najmniej n razy
<code>{n,m}</code>	Co najmniej n razy ale nie więcej niż m razy

Poniżej przedstawiono kilka innych elementów.

<code>\</code>	Ukośnik jest również znakiem ucieczki, tzn. likwiduje znaczenie znaków specjalnych, np. <code>\.</code> Znaczy dosłownie „kropka”
<code>\0x20</code>	Kod szesnastkowy znaku ASCII
<code>\u0020</code>	Kod szesnastkowy znaku Unicode
<code>()</code>	Nawiasy zwykle używane są do grupowania symboli
<code> </code>	Pionowa kreska oznacza alternatywę

**Uwaga:** Wyrażenia regularne są wpisywane jako łańcuchy znakowe (String), w których znak ukośnika (`\`) ma również specjalne znaczenie, dlatego wpisując wyrażenie regularne należy podwójnie wpisywać znak ukośnika jeśli ma on pełnić funkcje znaku ucieczki. Tzn. zamiast wpisać „`\`” wpisujemy „`\\`”, ale jeżeli znak ukośnika nie pełni funkcji znaku ucieczki to wpisujemy normalnie np. „`\d`”.

Przykład działania wyrażeń regularnych:

- sprawdzanie kodu pocztowego:  $^{\wedge}\backslash d\{2\}-\backslash d\{3\}\$$
- sprawdzenie adresu www:  $^{\wedge}www\backslash \cdot [-\backslash w]^+(\backslash \cdot [-\backslash w]^+)+\$$
- sprawdzanie adresu IP:  $^{\wedge}0^*(25[0-5]|2[0-4]\backslash d|1?\backslash d\backslash d?)(\backslash \cdot 0^*(25[0-5]|2[0-4]\backslash d|1?\backslash d\backslash d?))\{3\}\$$

Jeżeli chcemy dokonać dokładnego sprawdzenia zgodności wpisanego tekstu z wyrażeniem regularnym to na początku wyrażenia regularnego musimy umieścić znak  $^{\wedge}$  a na jego końcu znak  $\$$ .

Więcej informacji o wyrażeniach regularnych:

<http://www.codeguru.pl/article-696.aspx>

<http://www.radsoftware.com.au/articles/regexlearnsyntax.aspx>

## 2. Jak przejść z definicji gramatyki na wyrażenia regularne

Do przejścia na wyrażenia regularne najlepiej nadają się gramatyki zapisane w notacji MBNF w której dochodzą dodatkowe symbole:

[ ]	Oznacza wystąpienie jakiegoś symbolu zero lub jeden raz
{ }	Oznacza wystąpienie jakiegoś symbolu nieokreśloną ilość razy
( )	Oznacza grupę symboli

Przejście z gramatyki na wyrażenie regularne zostanie przedstawione na przykładzie prostej gramatyki dla adresu e-mail.

```
S ::= A@A.W
A ::= W{.W}
W ::= L{L}
L ::= a|b|c|d|e|...|x|y|z
```

Następnie należy kolejno po podstawiać do poszczególnych symboli. W rezultacie otrzymamy:

```
S ::= (a|b|c|...|z) { (a|b|c|...|z) } { . (a|b|c|...|z) { (a|b|c|...|z) } }
    @ (a|b|c|...|z) { (a|b|c|...|z) } { . (a|b|c|...|z) { (a|b|c|...|z) } }
    . (a|b|c|...|z) { (a|b|c|...|z) }
```

Z kolei ten zapis jest już bez mała wyrażeniem regularnym:

$^[a-z]+(\.[a-z]+)^*@[a-z]+(\.[a-z]+)^*\.[a-z]+\$$

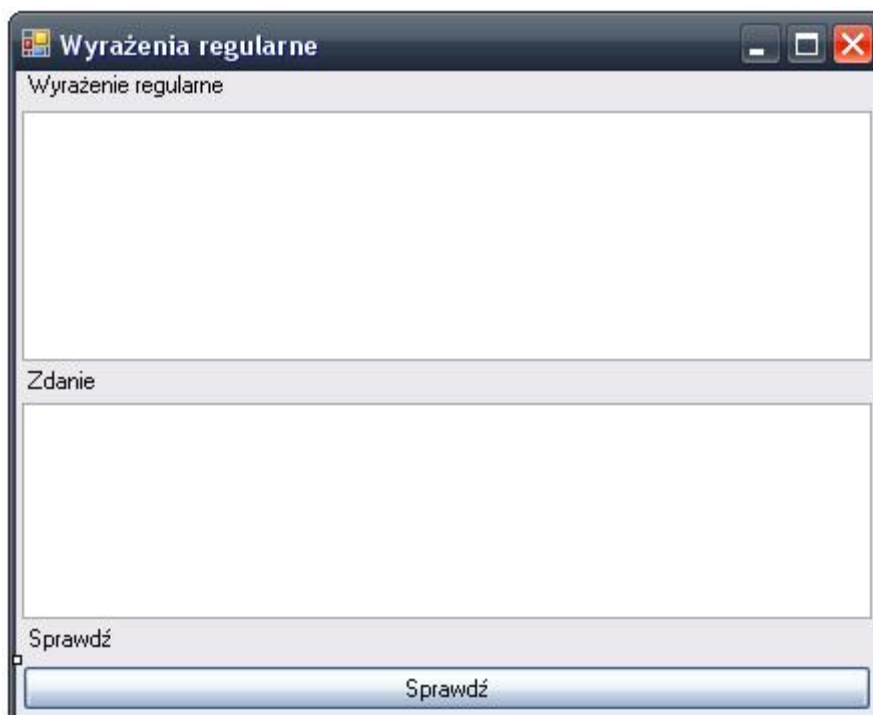
Niestety jeżeli w gramatyce występuje rekurencyjne wywołanie jakiegoś symbolu, to nie można dokonać podstawienia i nie da się zbudować wyrażenia regularnego do sprawdzenia tej gramatyki. Np.

```
S ::= W; {W;}
W ::= POP{OP}
P ::= L|(W)
L ::= C{C}
C ::= 1|2|3|4|5|6|7|8|9|0
O ::= *|/|+|-|^
```

Gdy podstawimy P do W to w W będziemy mieli rekurencyjne odwołanie do W. i nie będzie możliwe dokonanie dalszych podstawień.

### 3. Program

W trakcie laboratorium można korzystać z przygotowanego programu w którego jednym polu należy wpisać wyrażenie regularne a w drugim zdanie do sprawdzenia przy pomocy tego wyrażenia regularnego.



#### 4. Zadania

Napisać wyrażenia regularne dla sprawdzania zgodności zdań z zadanymi gramatykami:

$S ::= W; \{W\}$   
 $W ::= LOL\{OL\}$   
 $L ::= C\{C\}$   $12+2; 3*8^{12}-2/3;$   
 $C ::= 1|2|3|4|5|6|7|8|9|0$   
 $O ::= *|/|+|-|^$

$S ::= A[L][eA]$   
 $A ::= [-]L$   $-15.6e-3$   
 $L ::= C\{C\}$   
 $C ::= 1|2|3|4|5|6|7|8|9|0$

Wymyślić gramatykę i napisać dla niej wyrażenie regularne.