

Регулярные выражения

1

**Романенко Владимир Васильевич,
к.т.н., доцент каф. АСУ ТУСУР**

Способы определения языка

2

Способ 2. Непосредственная запись РВ. В этом случае можно действовать следующим образом:

Вариант 1 (последовательное построение РВ).

- Определить начальную конструкцию, которая может следовать в начале правильной цепочки в языке L . Составить РВ, соответствующее этой конструкции.
- Составить РВ для следующей конструкции, и т.д. до тех пор, пока не получим РВ, полностью описывающее язык L .

Способы определения языка

3

При этом можно использовать следующие типовые приёмы:

Приём 1. Конкатенация. Если конструкции a и b должны следовать друг за другом, то этому будет соответствовать часть РВ

$...(ab)...$

Приём 2. Объединение. Если следующей конструкцией входной цепочки является либо a , либо b , то этому будет соответствовать часть РВ

$...(a + b)...$

Способы определения языка

4

При этом можно использовать следующие типовые приёмы:

Приём 3. Положительная итерация. Если конструкция a в данной части входной цепочки может повторяться от 1 и более раз, то этому будет соответствовать часть PV

$$\dots(aa^*)\dots \text{ или } \dots(a^+)\dots$$

Приём 4. Обычная итерация. Если в предыдущем случае цепочка a может вообще отсутствовать, то

$$\dots(a^*)\dots$$

Способы определения языка

5

При этом можно использовать следующие типовые приёмы:

Приём 5. Конкатенация с итерацией. Если в данной части входной цепочки после конструкции a может повторяться от 0 и более раз другая конструкция b , то

$$\dots(ab^*)\dots$$

Если конструкция b должна повторяться от 1 и более раз, то

$$\dots(abb^*)\dots \text{ или } \dots(ab^+)\dots$$

Способы определения языка

6

При этом можно использовать следующие типовые приёмы:

Приём 6. Итерация объединения. Если в данной части входной цепочки может в произвольном порядке располагаться произвольное количество (от 0 и более) конструкций a и b , то

$$\dots(a+b)^*\dots$$

Если должна быть как минимум одна такая конструкция, то

$$\dots(a+b)(a+b)^*\dots \text{ или } \dots(a+b)^+\dots$$

Способы определения языка

7

При этом можно использовать следующие типовые приёмы:

Приём 7. Итерация конкатенации. Если в данной части входной цепочки может располагаться произвольное количество (от 0 и более) конструкций a и b , следующих друг за другом, то

$$\dots(ab)^*\dots$$

Если должно быть как минимум одно вхождение таких конструкций, то

$$\dots(ab)(ab)^*\dots \text{ или } \dots(ab)^+\dots$$

Примеры определения языка

8

Пример 1. В языке L , описывающем число с фиксированной точкой, в начале входной цепочки может располагаться:

- точка;
- **или** цифра;
- **или** знак (плюс или минус).

Поэтому получим следующее РВ:

$$.x + (0-9)y + (+ -)z$$

Т.е. (точка **и** x) **или** (цифра **и** y) **или** (знак (плюс **или** минус) **и** z).

Примеры определения языка

9

$$.x + (0-9)y + (+ + -)z$$

1. Начнём с точки (.). После точки могут следовать только цифры, от 1 и более, поэтому

$$x = (0-9)^+ \Rightarrow .x = .(0-9)^+$$

2. Далее, цифра (0-9). Это значит, у числа есть целая часть. После первой цифры могут следовать от 0 и более других цифр (цифра **и** от 0 и более других цифр), т.е. получим

$$y = (0-9)^*y_2 \Rightarrow (0-9)y = (0-9)(0-9)^*y_2 = (0-9)^+y_2$$

Примеры определения языка

10

$$.(0-9)^+ + (0-9)^+y_2 + (+ + -)z$$

После целой части может следовать:

- точка.
- **или** ничего (конец входной цепочки).

$$(0-9)^+y_2 = (0-9)^+.y_3 + (0-9)^+e$$

После точки могут следовать:

- цифры, от 1 и более (дробная часть числа);
- **или** ничего.

$$(0-9)^+.y_3 = (0-9)^+.(0-9)^+ + (0-9)^+.e$$

Примеры определения языка

11

$$.(0-9)^+ + (0-9)^+.(0-9)^+ + (0-9)^+ + (0-9)^+ + (+ + -)z$$

3. И последнее – знак, плюс или минус (+ + –).
После него может следовать точка или цифра и
т.д., т.е. уже полученные ранее выражения:

$$\begin{aligned} z &= .x + (0-9)y = \\ &= .(0-9)^+ + (0-9)^+.(0-9)^+ + (0-9)^+ + (0-9)^+ \end{aligned}$$

Окончательно имеем:

$$\begin{aligned} &.(0-9)^+ + (0-9)^+.(0-9)^+ + (0-9)^+ + (0-9)^+ + \\ &(+ + -)(.(0-9)^+ + (0-9)^+.(0-9)^+ + (0-9)^+ + (0-9)^+) \end{aligned}$$

Способы определения языка

12

Вариант 2 (построение РВ методом декомпозиции).

- Определяем конструкции верхнего уровня, т.е. самые главные блоки, из которых состоят цепочки языка.
- Каждый полученный блок декомпозируем на более мелкие блоки до тех пор, пока все блоки не будут записаны в виде элементарных операций РВ.

Примеры определения языка

13

Пример 2. Цепочки языка L , описывающего число с фиксированной точкой, на самом верхнем уровне состоят из необязательного знака и мантиссы:

[знак]мантисса

Т.е. знак может отсутствовать, либо это может быть плюс или минус:

(+ + - + e)мантисса

Примеры определения языка

14

Как мы уже выясняли ранее, мантисса может принимать четыре разные формы: « $N.M$ », « $N.$ », « $.M$ », « N », где N – целая, а M – дробная часть числа, т.е. получим

$$\text{мантисса} = (N.M + N. + .M + N)$$

Также мы отмечали, что числа N и M имеют одинаковый формат – это последовательность из одной и более цифр в диапазоне от 0 до 9, т.е.

$$N = M = (0-9)^+.$$

Примеры определения языка

15

Окончательно получим

$$\begin{aligned} & (+ \text{ } + \text{ } - \text{ } + \text{ } e)(N.M \text{ } + \text{ } N. \text{ } + \text{ } .M \text{ } + \text{ } N) = \\ & = (+ \text{ } + \text{ } - \text{ } + \text{ } e)((0-9)^+.(0-9)^+ \text{ } + \text{ } (0-9)^+. \text{ } + \text{ } .(0-9)^+ \text{ } + \text{ } (0-9)^+) \end{aligned}$$

Оба полученных выражения можно упростить, чтобы они совпали с полученным способом №1 выражением

$$(+ \text{ } + \text{ } - \text{ } + \text{ } e)(.(0-9)^+ \text{ } + \text{ } (0-9)^+(e \text{ } + \text{ } .(0-9)^*))$$

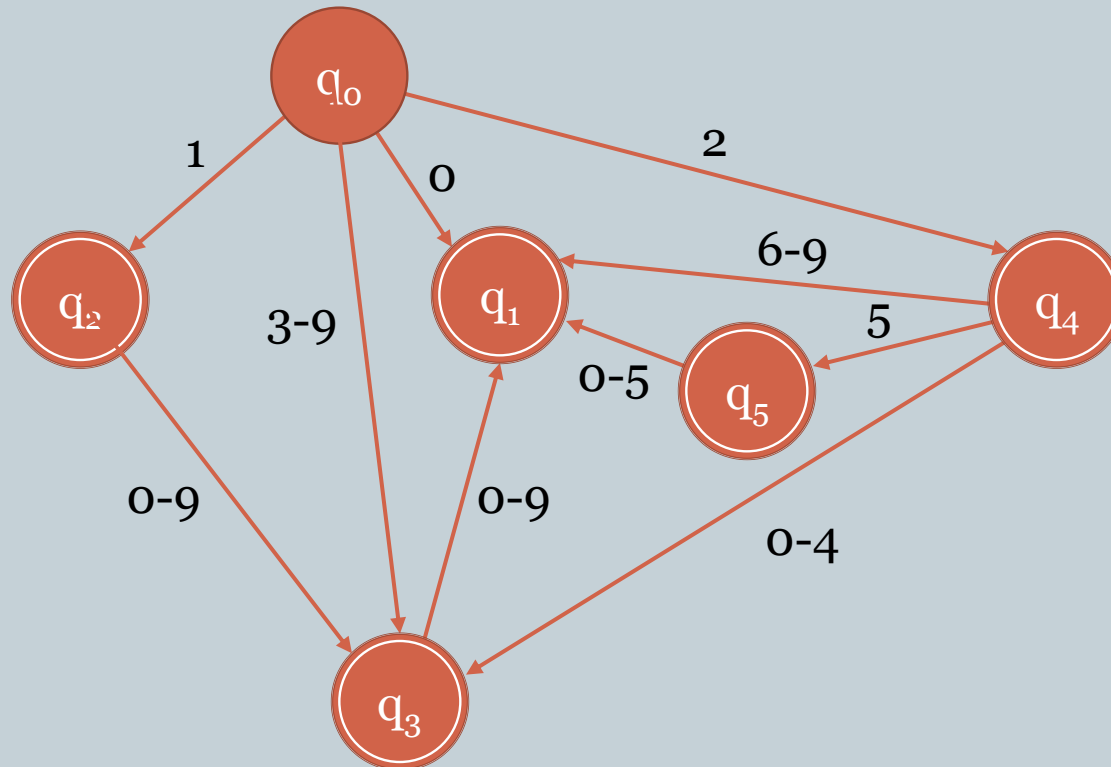
При этом полагается, что

$$0-9 = 0 \text{ } + \text{ } 1 \text{ } + \text{ } 2 \text{ } + \text{ } 3 \text{ } + \text{ } 4 \text{ } + \text{ } 5 \text{ } + \text{ } 6 \text{ } + \text{ } 7 \text{ } + \text{ } 8 \text{ } + \text{ } 9$$

Примеры определения языка

16

Пример 3. Язык L описывает десятичные числа в диапазоне от 0 до 255, без ведущих нулей.



Примеры определения языка

17

Получим РВ:

$$\begin{aligned} & 0 + 1(e + (0-9) + (0-9)(0-9)) + \\ & + 2(e + (0-4)(e + (0-9)) + 5(e + (0-5)) + (6-9)) + \\ & + (3-9)(e + (0-9)) \end{aligned}$$

Примеры определения языка

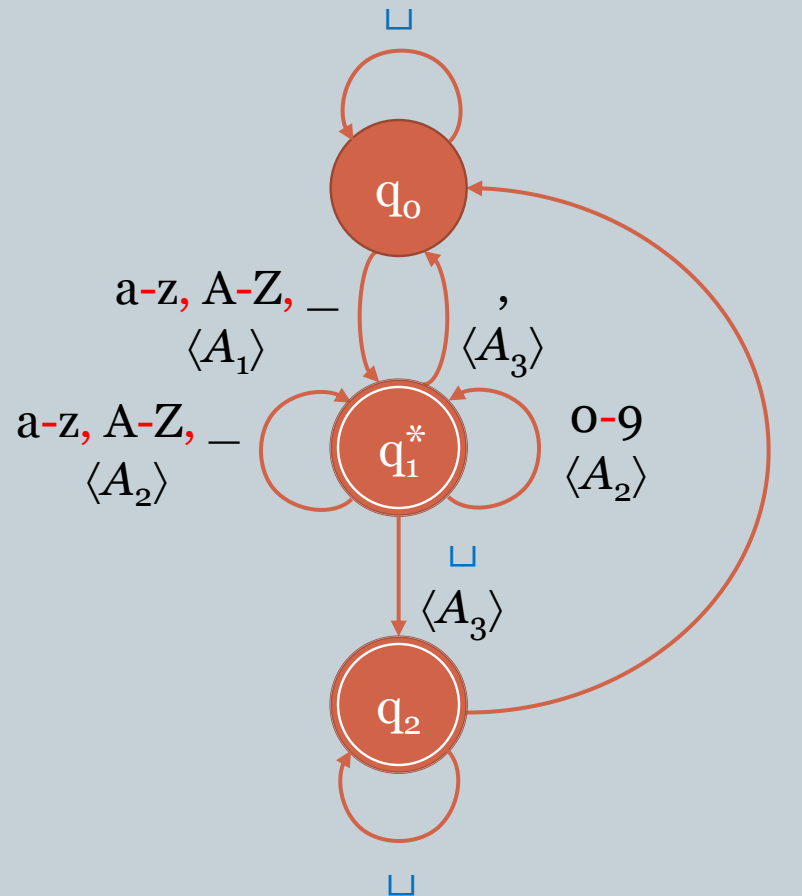
18

Пример 4. Пусть язык L описывает список идентификаторов, разделенных запятыми. Идентификатор должен начинаться с буквы латинского алфавита или подчёркивания, далее могут следовать другие необязательные буквы, подчёркивания или цифры. Также в этой записи допустимы пробелы и другие символы-разделители в начале и конце входной цепочки, а также между идентификаторами и запятыми. Идентификаторы должны быть уникальными.

Примеры определения языка

19

Граф функции переходов:



Примеры определения языка

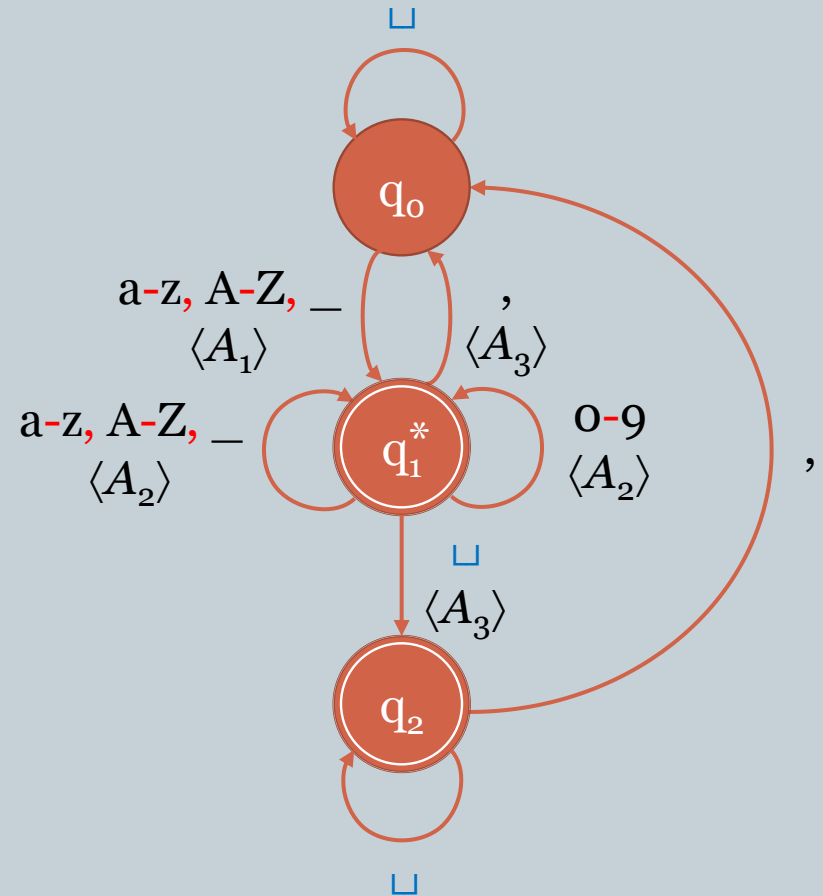
20

id, id, ..., id =
 = (id,)*id или id(id,)*

$\square^* \text{id} \square^*, \square^* \text{id} \square^*, \square^* \dots$
 $\square^*, \square^* \text{id} \square^* =$
 $(\square^* \text{id} \square^*, \square^*)^* \square^* \text{id} \square^* =$
 $(\square^* \text{id} \square^*,)^* \square^* \text{id} \square^*$

или

$\square^* \text{id} \square^* (, \square^* \text{id} \square^*)^*$



Примеры определения языка

21

Получим РВ:

$$\begin{aligned} & \square^* \text{id } \square^* (, \square^* \text{id } \square^*)^* = \\ & = \square^* (_ + a\text{-}z + A\text{-}Z)(_ + a\text{-}z + A\text{-}Z + 0\text{-}9)^* \square^* \\ & (, \square^* (_ + a\text{-}z + A\text{-}Z)(_ + a\text{-}z + A\text{-}Z + 0\text{-}9)^* \square^*)^* \end{aligned}$$

При этом полагается, что

$$0\text{-}9 = 0 + 1 + 2 + 3 + 4 + 5 + 6 + 7 + 8 + 9,$$

$$a\text{-}z = a + b + c + \dots + z,$$

$$A\text{-}Z = A + B + C + \dots + Z$$

Программирование РВ

22

Все современные языки программирования позволяют работать с регулярными выражениями. При этом часть языков (PHP, JavaScript и т.д.) имеют встроенные средства для работы с РВ, а в других есть готовые библиотеки или классы для работы с ними, например, класс `Regex` для платформы .NET:

```
using System.Text.RegularExpressions;  
Regex r = new Regex("шаблон");
```

Программирование РВ

23

Символ	Интерпретация
Escape-последовательности	
<code>\b</code>	При использовании его в квадратных скобках соответствует символу «←» (<code>\u0008</code>)
<code>\t</code> , <code>\r</code> , <code>\n</code> , <code>\a</code> , <code>\f</code> , <code>\v</code>	Табуляция (<code>\u0009</code>), возврат каретки (<code>\u000D</code>), новая строка (<code>\u000A</code>) и т.д.
<code>\cX</code>	Управляющий символ (например, <code>\cC</code> – это Ctrl+C, <code>\u0003</code>)
<code>\e</code>	Escape (<code>\u001B</code>)
<code>\ooo</code>	Символ ASCII в восьмеричной системе
<code>\xhh</code>	Символ ASCII в шестнадцатеричной системе
<code>\uhhhh</code>	Символ Unicode
<code>\</code>	Следующий символ не является специальным символом РВ. Этим символом нужно экранировать все специальные символы

Пример (в примере приведен шаблон и строка поиска, в строке найденные совпадения подчеркнуты): `@\"\\r\\n\\w+\" – \"\\r\\nЗдесь имеются\\nдве строки\".`

Программирование РВ

24

Подмножества символов

.	Любой символ, кроме конца строки (\n)
[xxx]	Любой символ из множества
[^xxx]	Любой символ, кроме символов из множества
[x-x]	Любой символ из диапазона
[xxx-[xxx]]	Вычитание одного множества или диапазона из другого
\p{name}	Любой символ, заданный категорией Unicode с именем name
\P{name}	Любой символ, кроме заданных категорией Unicode с именем name
\w	Множество символов, используемых при задании идентификаторов
\W	Множество символов, не используемых при задании идентификаторов
\s	Пробелы
\S	Все, кроме пробелов
\d	Цифры
\D	Не цифры

Примеры: @".+" – "\r\nЗдесь имеются\ndве строки"; @[fx]+" – "oxabcfx"; @"[^fx]+" – "oxabcfx";
@"[a-f]+" – "oxabcfx"; @"[^a-f]+" – "oxabcfx"; @"[a-z-[c]]+" – "oxabcfx"; @"\p{Lu}" – "City Lights";
@"\P{Lu}" – "City"; @"\p{IsCyrillic}" – "хаОС"; @"\P{IsCyrillic}" – "хаОС". // Lu – прописные
латинские буквы, IsCyrillic – русские буквы

Программирование РВ

25

Привязка	
^, \A	В начале строки
\$, \Z	В конце строки или до символа «\n» в конце строки
\z	В конце строки
\G	В том месте, где заканчивается предыдущее соответствие
\b	Граница слова
\B	Любая позиция не на границе слова
Примеры:	
@"\G\(\\d\\)" – "(1)(3)(5)[7](9) "; // три соответствия (1), (2) и (3)	
@"\bn\\S*ion\\b" – "<u>nation</u> donation";	
@"\Bend\\w*\\b" – "end <u>sends</u> endure <u>lender</u>".	

Программирование РВ

26

Операции (кванторы)

*, *?	Итерация
+, +?	Положительная итерация
?, ??	Ноль или одно соответствие
{n}, {n}?	Точно n соответствий
{n,}, {n,}?	По меньшей мере, n соответствий
{n,m},{n,m}?	От n до m соответствий

Примеры (первые кванторы – жадные, ищут как можно большее число элементов, вторые – ленивые, ищут как можно меньшее число элементов):

@"\d{3,}" – "888-555-5555";

@"^d{3}" – "913-913-913";

@"-d{3}\$" – "913-913-913";

@"5+?5" – "888-555-5555"; // три совпадения – 55, 55 и 55

@"5+5" – "888-555-5555".

Программирование РВ

27

Группирование	
()	Группа, автоматически получающая номер
(?:)	Не сохранять группу
(?<имя>) или (?'имя')	При обнаружении соответствия создается именованная группа
(?<имя–имя>) или (?'имя–имя')	Удаление ранее определенной группы и сохранение в новой группе подстроки между ранее определенной группой и новой группой
(?imnsx:) (?–imnsx:)	Включает или выключает в группе любую из пяти возможных опций: i – нечувствительность к регистру; s – одна строка (тогда «.» – это любой символ); m – многострочный режим («^», «\$» – начало и конец каждой строки); n – не захватывать неименованные группы; x – исключить не преобразованные в escape-последовательность пробелы из шаблона и включить комментарии после знака номера (#)
(?=)	Положительное утверждение просмотра вперед нулевой длины

Программирование РВ

28

Группирование (продолжение)

(?!)	Отрицательное утверждение просмотра вперед нулевой длины
(?<=)	Положительное утверждение просмотра назад нулевой длины
(?<!)	Отрицательное утверждение просмотра назад нулевой длины
(?>)	Невозвращаемая (жадная) часть выражения

Примеры:

@"(an)+" – "bananas annals";

@"an+" – "bananas annals"; // сравните, три совпадения – an, an и ann

@"(?i:an)+" – "baNAnas annals";

@"[a-z]+(?=\d)" – "abc xyz12 555w";

@"(?<=\d)[a-z]+" – "abc xyz12 555w".

Программирование РВ

29

Ссылки

`\число`

Ссылка на группу

`\k<имя>`

Ссылка на именованную группу

Примеры:

`@"(\w)\1" – "deep";`

`@"(?<char>\w)\k<char>" – "deep".`

Конструкции изменения

`|`

Альтернатива (соответствует операции объединения)

`(?(выражение)да|нет)`

Сопоставляется с частью «да», если выражение соответствует; в противном случае сопоставляется с необязательной частью «нет»

`(?(имя)да|нет),`

`(?(число)да|нет)`

Сопоставляется с частью «да», если названное имя захвата имеет соответствие; в противном случае сопоставляется с необязательной частью «нет»

Пример:

`@"th(e|is|at)" – "this is the day";`

Программирование РВ

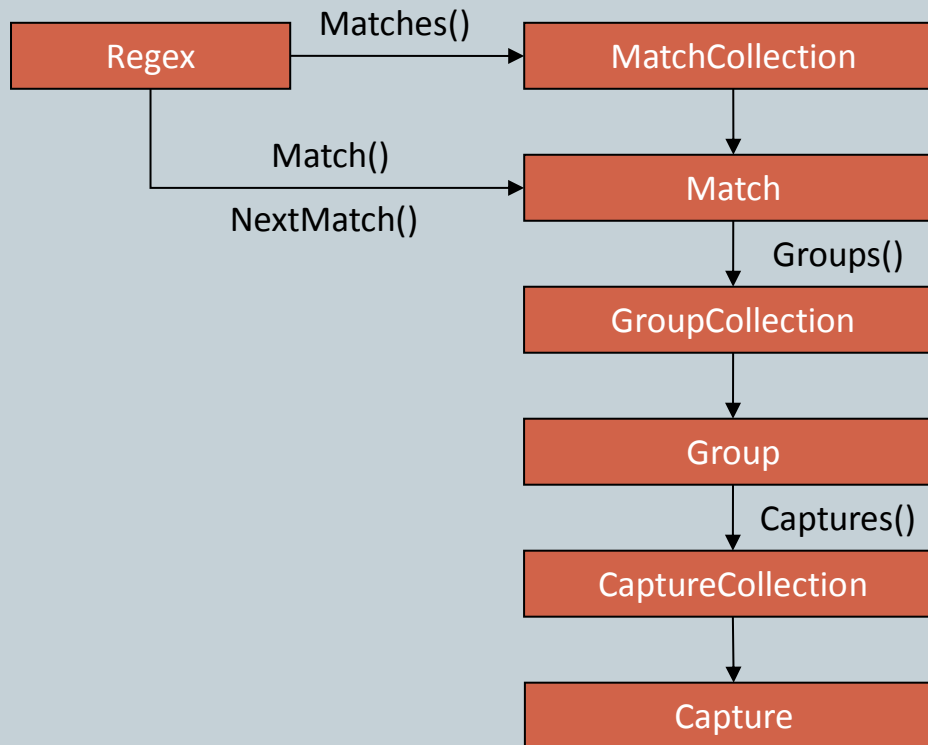
30

Подстановки	
\$число	Замещается часть строки, соответствующая группе с указанным номером
\${имя}	Замещается часть строки, соответствующая группе с указанным именем
\$\$	Подставляется \$
\$&	Замещение копией полного соответствия
\$`	Замещение текста входной строки до соответствия
\$'	Замещение текста входной строки после соответствия
\$+	Замещение последней захваченной группы
\$_	Замещение всей строки
Комментарии	
(?#)	Встроенный комментарий
#	Комментарий до конца строки

Программирование РВ

31

Результаты работы Regex:



Программирование РВ

32

Пример:

```
Regex r = new Regex(@"((\d)+)+");
Match m = r.Match("123 456");
int matchCount = 0;
while (m.Success)
{
    Console.WriteLine("Соответствие {0}", ++matchCount);
    for (int i = 1; i < m.Groups.Count; i++)
    {
        Group g = m.Groups[i];
        Console.WriteLine("    Группа {0} = '{1}'", i, g.Value);
        for (int j = 0; j < g.Captures.Count; j++)
        {
            Capture c = g.Captures[j];
            Console.WriteLine("        Захват {0} = '{1}', позиция = {2},
        длина = {3}", j, c, c.Index, c.Length);
        }
        m = m.NextMatch();
    }
}
```

Соответствие 1

Группа 1 = '123'

Захват 0 = '123', позиция = 0, длина = 3

Группа 2 = '3'

Захват 0 = '1', позиция = 0, длина = 1

Захват 1 = '2', позиция = 1, длина = 1

Захват 2 = '3', позиция = 2, длина = 1

Соответствие 2

Группа 1 = '456'

Захват 0 = '456', позиция = 4, длина = 3

Группа 2 = '6'

Захват 0 = '4', позиция = 4, длина = 1

Захват 1 = '5', позиция = 5, длина = 1

Захват 2 = '6', позиция = 6, длина = 1

Программирование РВ

33

Пример 1. Язык L , описывающий двоичные числа без незначащих нулей:

$$0 + 1(0 + 1)^* \Rightarrow "\textcolor{red}{^}(0|1(0|1)^*)\textcolor{red}{\$}"$$

Программирование РВ

34

Пример 2. Язык L , описывающий числа с фиксированной точкой:

$$(+ \text{ } + \text{ } - \text{ } + \text{ } e)(.(0-9)^+ + (0-9)^+(e \text{ } + \text{ } .(0-9)^*))$$

Во-первых, $0-9 \Rightarrow 0|1|2|3|4|5|6|7|8|9 = [0-9] = \backslash d$

Во-вторых, $+ \text{ } + \text{ } - \text{ } + \text{ } e \Rightarrow (\backslash + | -)?$

$$e \text{ } + \text{ } .(0-9)^* \Rightarrow (\backslash . \backslash d^*)?$$

Таким образом,

$$\Rightarrow "\text{ }^(\backslash \backslash + | -)?(\backslash \backslash . \backslash \backslash d^+ | \backslash \backslash d^+(\backslash \backslash . \backslash \backslash d^*)?)\$"$$
 или
$$@"\text{ }^(\backslash + | -)?(\backslash . \backslash d^+ | \backslash d^+(\backslash . \backslash d^*)?)\$"$$

Программирование РВ

35

Пример 3. Язык L описывает десятичные числа в диапазоне от 0 до 255, без ведущих нулей:

$$\begin{aligned} & 0 + 1(e + (0-9) + (0-9)(0-9)) + \\ & + 2(e + (0-4)(e + (0-9)) + 5(e + (0-5)) + (6-9)) + \\ & + (3-9)(e + (0-9)) \Rightarrow \end{aligned}$$

$$\begin{aligned} @'^{(0|1\d{0,2}|2([0-4]\d?|5[0-5]?|[6-9])?| \\ [3-9]\d?)\$"} \end{aligned}$$

Или

$$\d{0,2} = \d?\d?$$

Программирование РВ

36

Пример 4. Язык L описывает список **уникальных** идентификаторов:

$$\begin{aligned} & \square^* (_ + a-z + A-Z)(_ + a-z + A-Z + 0-9)^* \square^* \\ & (, \square^* (_ + a-z + A-Z)(_ + a-z + A-Z + 0-9)^* \square^*)^* \\ \Rightarrow & @'^{\wedge} \backslash s^* [_ a-zA-Z][_ a-zA-Z0-9]^* \backslash s^* \\ & (, \backslash s^* [_ a-zA-Z][_ a-zA-Z0-9]^* \backslash s^*)^* \$'' \end{aligned}$$

или

$$\begin{aligned} & @'^{\wedge} (((?<!^), | (?<=^)) \backslash s^* \\ & [_ a-zA-Z][_ a-zA-Z\d]^* \backslash s^*)^* \$'' \end{aligned}$$

Программирование РВ

37

Пример 4. Язык L описывает список **уникальных** идентификаторов:

$$\begin{aligned} & \square^* (_ + a-z + A-Z)(_ + a-z + A-Z + 0-9)^* \square^* \\ & (, \square^* (_ + a-z + A-Z)(_ + a-z + A-Z + 0-9)^* \square^*)^* \\ \Rightarrow & @'^\wedge \backslash s^* [_ a-zA-Z][_ a-zA-Z0-9]^* \backslash s^* \\ & (, \backslash s^* [_ a-zA-Z][_ a-zA-Z0-9]^* \backslash s^*)^* \$'' \end{aligned}$$

или (с русскими буквами)

$$@'^\wedge \backslash s^* [_ \backslash p\{L\}] \backslash w^* \backslash s^* (, \backslash s^* [_ \backslash p\{L\}] \backslash w^* \backslash s^*)^* \$''$$

Внедрение действий в синтаксис

38

Внедрение действий осуществляется с помощью анализа групп вхождений:

- нумерованных;
- именованных.

С нумерованными группами работать сложнее.

Внедрение действий в синтаксис

39

Пример 1. Язык L описывает список **уникальных** идентификаторов:

$$@'^{(((?<!^),|(?<=^))\backslash s^*[_a-zA-Z][_a-zA-Z\backslash d]^*\backslash s^*)^*\$"$$

Входная строка «L1b2 , cdX34»:

Соответствие 1

Группа 1 = ' , cdX34 '

Захват 0 = ' L1b2 ', позиция = 0, длина = 6

Захват 1 = ' , cdX34 ', позиция = 6, длина = 8

Группа 2 = ' , '

Захват 0 = ' ', позиция = 0, длина = 0

Захват 1 = ' , ', позиция = 6, длина = 1

Внедрение действий в синтаксис

40

Модифицируем выражение:

$$@'^{((?<!^),|(?<=^))\backslash s^*}$$
$$([_a-zA-Z][_a-zA-Z\backslash d]^*)\backslash s^*)^*\$"$$

Соответствие 1

Группа 1 = ' cdX34 '

Захват 0 = ' L1b2 ', позиция = 0, длина = 6

Захват 1 = ' cdX34 ', позиция = 6, длина = 8

Группа 2 = ' , '

Захват 0 = ' ', позиция = 0, длина = 0

Захват 1 = ' , ', позиция = 6, длина = 1

Группа 3 = ' cdX34 '

Захват 0 = ' L1b2 ', позиция = 1, длина = 4

Захват 1 = ' cdX34 ', позиция = 8, длина = 5

Внедрение действий в синтаксис

41

Введём именованную группу:

`@'^(((?<!^),|(?<=^))\s*
(?id>[_a-zA-Z][_a-zA-Z\d]*)\s*)*$"`

```
Match m = r.Match(" L1b2 , cdX34 ");  
foreach (Capture c in m.Groups["id"].Captures)  
{  
    Console.WriteLine(c.Value);  
}
```

L1b2
cdX34

Внедрение действий в синтаксис

42

Пример 2. Язык L описывает десятичные числа в диапазоне от 0 до 255, без ведущих нулей:

`@"^(0|1\d{0,2}|2([0-4]\d?|5[0-5]?|[6-9])?|[3-9]\d?)$"`

Введём именованную группу:

`@"^(?num\d+)$"`

```
Match m = r.Match("123");
int matchCount = 0;
int value;
if (!int.TryParse(m.Groups["num"].Value, out value) ||
    value > 255) Console.WriteLine("Ошибка!");
```

Сбалансированные определения

Для проверки рекурсивно вложенных описаний (таких, для анализа которых используется ДМПА) применяются *сбалансированные определения*. В этом случае именованные группы играют роль стека. Так, конструкция «('x')» добавляет в коллекцию с именем «x» один элемент, а конструкция «('x')» убирает из коллекции «x» один элемент. Поэтому в конце остается лишь проверить, что в коллекции не осталось элементов – «(? (x) (!))».

Сбалансированные определения

44

Пример. Пусть язык L описывает вложенные операторы языка Pascal «**begin end**;». Учитывая, что необходимо проверять их парность, используем ДМПА с посимвольным разбором. Операторы отделяются друг от друга разделительными символами (пробелами, табуляциями, знаками возврата каретки и перехода на новую строку) в произвольном количестве, но не менее одного (в таблице обозначены символом подчеркивания). Также пробелы могут окружать знак «;».

Сбалансированные определения

45

Получим следующее РВ:

```
@'^\s*((?'begin'begin\s+)+  
(?'-begin'end\s*;\s*)+)*(?begin)(?!))$"
```

Определение позиции ошибки

46

Структура РВ должна быть такой, чтобы не просто проверить соответствие шаблона входной цепочке, но найти максимальное количество соответствий, пусть даже частичных. Тогда позицией ошибки будет:

- 1) первая позиция входной цепочки (1), если первое соответствие не начинается с позиции $\text{Index} = 0$;
- 2) позиция, следующая за последним соответствием ($\text{Index} + \text{Length} + 1$), если она не совпадает с последней позицией входной цепочки;
- 3) позиция первого разрыва между соответствиями, если символ, следующий за предыдущим соответствием, не является первым символом следующего соответствия;
- 4) позиция специального захвата некорректных цепочек.

Определение позиции ошибки

47

Пример.

```
Regex r = new Regex(@"\w+(\.\w+)*");
string str = "abc.xyz.pqr";
MatchCollection m = r.Matches(str);
if (m.Count == 1 && m[0].Value == str) Console.WriteLine("OK");
else if (m.Count == 0) Console.WriteLine("Ошибка в позиции 1 '{0}'",
str[0]);
else
{
    int index = 0;
    for (int i = 0; i < m.Count; i++)
    {
        if (m[i].Index > index) break;
        index = m[i].Index + m[i].Length;
    }
    Console.WriteLine("Ошибка в позиции {0} '{1}'", index + 1,
str[index]);
}
```

Определение позиции ошибки

48

Результаты работы:

- «abc.xyz.pqr» – правильно;
- «+abc.xyz.pqr» – ошибка в позиции 1 («+»);
- «abc.xyz.pqr!» – ошибка в позиции 12 («!»);
- «abc.xyz!.pqr» – ошибка в позиции 8 («!»).
- «abc.xyz.+pqr» – **ошибка в позиции 8 («.»).**

Определение позиции ошибки

49

Исправляем:

@\"\\w+(\\.\\w+)*\\.(?!\$))?"

Результаты работы:

- «abc.xyz.+pqr» – ошибка в позиции 9 («+»);
- «abc.xyz.pqr.» – ошибка в позиции 12 («.»).

Лабораторная работа №2

50

Порядок выполнения лабораторной работы:

1. Описать требуемый язык заданным способом (в виде РВ).
2. Написать программу, реализующую требуемый механизм синтаксического анализа.
3. Внедрить в синтаксис анализатора действия для проверки семантики языка или его интерпретации.
4. Протестировать программу.
5. Написать отчёт, включающий все требуемые пункты (в т.ч. формальное описание построенного анализатора) и удовлетворяющий требованиям ОС ТУСУР 01-2013.

Лабораторная работа №2

51

Требования к программе:

- В программе должно быть описано **единственное** РВ, полностью описывающее заданный язык.
- Для анализа РВ рекомендуется применять уже готовые классы и библиотеки (**Regex** и т.п.).
- Внедрение в синтаксис РВ действий и поиск позиции ошибки обеспечивается добавлением нумерованных и именованных групп:

(...) (**?<имя>**...) (**?'**имя**'**...)

- Требования к входным и выходным данным – лабораторной работе №1.

Лабораторная работа №1

52

Способы поиска ошибок при помощи РВ:

- Построить РВ так, чтобы выделить все правильные совпадения во входной цепочке и искать разрывы между ними:

```
int a, 12b, c3[10], x[];
```

- Построить РВ так, чтобы выделить правильную часть цепочки, начиная от её начала:

```
int a, 12b, c3[10], x[];
```

Для этого используются утверждения просмотра вперёд и назад нулевой длины:

```
(?!...) (?<=...) (?<!...) (?>...)
```

Домашнее задание

53

Язык L описывает обращение к элементу массива. При этом размерность массива может быть любой, а в качестве индексов можно использовать целые константы $i \geq 0$, а также идентификаторы, в т.ч. элементы других массивов. Идентификатор начинается с латинской буквы, после которой могут следовать другие буквы и цифры. Индексы заключаются в квадратные скобки и отделяются друг от друга запятыми. Например:

- `a[1];`
- `a2[5,b[2],z];`
- `mas[x[4],y[4]]` и т.д.

Домашнее задание

54

Вариант 1:

```
const string pattern = @"^
(
    (
        (?<=^|,|\[)      (?<op>[a-z]+)  (?=,|\[|\]|)  ) |
        (?<=,|\[)        (?<op>\d+)    (?=,|\]|)      ) |
        (?<=[a-z]|\d|\]|) (?<op>,)      (?=[a-z]|\d)    ) |
        (?<level>        (?<=[a-z])      (?<op>\[)        (?=[a-z]|\d)    ) |
        (?<-level>      (?<=[a-z]|\d|\]|) (?<op>\]|)        (?=,|\]||$)      )
    )+
    (? (level) (?!))
) $";

Regex r = new Regex(pattern, RegexOptions.IgnorePatternWhitespace);
```

Домашнее задание

55

Вариант 1:

```
abc[12,zz[10],tt,zx[zy[5],q]]
Соответствие 1
  Группа 1 = 'abc[12,zz[10],tt,zx[zy[5],q]]'
    Захват 0 = 'abc[12,zz[10],tt,zx[zy[5],q]]', позиция =
0, длина = 29
    Группа 2 = ']'
      Захват 0 = 'abc', позиция = 0, длина = 3
      Захват 1 = '[', позиция = 3, длина = 1
      Захват 2 = '12', позиция = 4, длина = 2
      Захват 3 = ',', позиция = 6, длина = 1
      Захват 4 = 'zz', позиция = 7, длина = 2
      Захват 5 = '[', позиция = 9, длина = 1
      Захват 6 = '10', позиция = 10, длина = 2
      Захват 7 = ']', позиция = 12, длина = 1
      Захват 8 = ',', позиция = 13, длина = 1
      Захват 9 = 'tt', позиция = 14, длина = 2
      Захват 10 = ',', позиция = 16, длина = 1
      Захват 11 = 'zx', позиция = 17, длина = 2
      Захват 12 = '[', позиция = 19, длина = 1
      Захват 13 = 'zy', позиция = 20, длина = 2
      Захват 14 = '[', позиция = 22, длина = 1
      Захват 15 = '5', позиция = 23, длина = 1
      Захват 16 = ']', позиция = 24, длина = 1
      Захват 17 = ',', позиция = 25, длина = 1
      Захват 18 = 'q', позиция = 26, длина = 1
      Захват 19 = ']', позиция = 27, длина = 1
      Захват 20 = ']', позиция = 28, длина = 1
    Группа 3 = 'q'
      Захват 0 = 'abc', позиция = 0, длина = 3
      Захват 1 = 'zz', позиция = 7, длина = 2
      Захват 2 = 'tt', позиция = 14, длина = 2
      Захват 3 = 'zx', позиция = 17, длина = 2
      Захват 4 = 'zy', позиция = 20, длина = 2
      Захват 5 = 'q', позиция = 26, длина = 1
```

```
Группа 4 = '5'
  Захват 0 = '12', позиция = 4, длина = 2
  Захват 1 = '10', позиция = 10, длина = 2
  Захват 2 = '5', позиция = 23, длина = 1
  Группа 5 = ','
    Захват 0 = ',', позиция = 6, длина = 1
    Захват 1 = ',', позиция = 13, длина = 1
    Захват 2 = ',', позиция = 16, длина = 1
    Захват 3 = ',', позиция = 25, длина = 1
    Группа 6 = ']'
      Захват 0 = 'abc', позиция = 0, длина = 3
      Захват 1 = '[', позиция = 3, длина = 1
      Захват 2 = '12', позиция = 4, длина = 2
      Захват 3 = ',', позиция = 6, длина = 1
      Захват 4 = 'zz', позиция = 7, длина = 2
      Захват 5 = '[', позиция = 9, длина = 1
      Захват 6 = '10', позиция = 10, длина = 2
      Захват 7 = ']', позиция = 12, длина = 1
      Захват 8 = ',', позиция = 13, длина = 1
      Захват 9 = 'tt', позиция = 14, длина = 2
      Захват 10 = ',', позиция = 16, длина = 1
      Захват 11 = 'zx', позиция = 17, длина = 2
      Захват 12 = '[', позиция = 19, длина = 1
      Захват 13 = 'zy', позиция = 20, длина = 2
      Захват 14 = '[', позиция = 22, длина = 1
      Захват 15 = '5', позиция = 23, длина = 1
      Захват 16 = ']', позиция = 24, длина = 1
      Захват 17 = ',', позиция = 25, длина = 1
      Захват 18 = 'q', позиция = 26, длина = 1
      Захват 19 = ']', позиция = 27, длина = 1
      Захват 20 = ']', позиция = 28, длина = 1
    Группа 7 = ''
```

Домашнее задание

56

Вариант 2:

```
const string pattern = @"^
    (?
        (
            (?
                (?<=^|,|\[)      (?<op>[a-z]+)  (?=,|\[|\])  ) |
                (?<=,|\[)      (?<op>\d+)      (?=,|\])      ) |
                (?<=[a-z]|\d|\]) (?<op>, )      (?=[a-z]|\d)  ) |
                (?<level>      (?<=[a-z])      (?<op>\[)      (?=[a-z]|\d)  ) |
                (?<-level>    (?<=[a-z]|\d|\]) (?<op>\])      (?=,|\]|\$)   )
            )+
        )?(level)?(!)
    )$";

Regex r = new Regex(pattern, RegexOptions.IgnorePatternWhitespace);
```


Домашнее задание

57

Вариант 2:

```
abc[12,zz[10],tt,zx[zy[5],q]]
```

Соответствие 1

Группа 1 = ']'

Захват 0 = 'abc', позиция = 0, длина = 3

Захват 1 = '[', позиция = 3, длина = 1

Захват 2 = '12', позиция = 4, длина = 2

Захват 3 = ',', позиция = 6, длина = 1

Захват 4 = 'zz', позиция = 7, длина = 2

Захват 5 = '[', позиция = 9, длина = 1

Захват 6 = '10', позиция = 10, длина = 2

Захват 7 = ']', позиция = 12, длина = 1

Захват 8 = ',', позиция = 13, длина = 1

Захват 9 = 'tt', позиция = 14, длина = 2

Захват 10 = ',', позиция = 16, длина = 1

Захват 11 = 'zx', позиция = 17, длина = 2

Захват 12 = '[', позиция = 19, длина = 1

Захват 13 = 'zy', позиция = 20, длина = 2

Захват 14 = '[', позиция = 22, длина = 1

Захват 15 = '5', позиция = 23, длина = 1

Захват 16 = ']', позиция = 24, длина = 1

Захват 17 = ',', позиция = 25, длина = 1

Захват 18 = 'q', позиция = 26, длина = 1

Захват 19 = ']', позиция = 27, длина = 1

Захват 20 = ']', позиция = 28, длина = 1

Группа 2 = ''