

## Лекция Регулярные выражения

Стандартный класс *string* позволяет выполнять над строками различные операции, в том числе поиск, замену, вставку и удаление подстрок. Тем не менее, есть классы задач по обработке символьной информации, где стандартных возможностей явно не хватает. Чтобы облегчить решение подобных задач, в *Net Framework* встроен более мощный аппарат работы со строками, основанный на регулярных выражениях.

Регулярные выражения предназначены для обработки текстовой информации и обеспечивают:

1. Эффективный поиск в тексте по заданному шаблону;
2. Редактирование текста;
3. Формирование итоговых отчетов по результатам работы с текстом.

Подробно рассмотрим первые два аспекта применения регулярных выражений.

### Метасимволы в регулярных выражениях

Регулярное выражение - это шаблон, по которому выполняется поиск соответствующего фрагмента текста. Язык описания регулярных выражений состоит из символов двух видов: обычных символов и метасимволов. Обычный символ представляет в выражении сам себя, а метасимвол - некоторый класс символов.

Рассмотрим наиболее употребительные метасимволы:

Класс символов	Описание	Пример
.	Любой символ, кроме <code>\n</code> .	Выражение <code>c.t</code> соответствует фрагментам: <code>cat</code> , <code>cut</code> , <code>c#t</code> , <code>c{t</code> и т.д.
[ ]	Любой одиночный символ из последовательности, записанной внутри скобок. Допускается использование диапазонов символов.	Выражение <code>c[ai]t</code> соответствует фрагментам: <code>cat</code> , <code>cut</code> , <code>cit</code> . Выражение <code>c[a-c]t</code> соответствует фрагментам: <code>cat</code> , <code>cbt</code> , <code>cct</code> .
[ ^ ]	Любой одиночный символ, не входящий в последовательность, записанную внутри скобок. Допускается использование диапазонов символов.	Выражение <code>c[^ai]t</code> соответствует фрагментам: <code>cbt</code> , <code>cct</code> , <code>c2t</code> и т.д. Выражение <code>c[^a-c]t</code> соответствует фрагментам: <code>cdt</code> , <code>cet</code> , <code>c%t</code> и т.д.
\w	Любой алфавитно - цифровой символ.	Выражение <code>c\wt</code> соответствует фрагментам: <code>cbt</code> , <code>cct</code> , <code>c2t</code> и т.д., но не соответствует фрагментам <code>c%t</code> , <code>c{t</code> и т.д.
\W	Любой не алфавитно - цифровой символ.	Выражение <code>c\Wt</code> соответствует фрагментам: <code>c%t</code> , <code>c{t</code> , <code>c.t</code> и т.д., но не соответствует фрагментам <code>cbt</code> , <code>cct</code> , <code>c2t</code> и т.д.
\s	Любой пробельный символ.	Выражение <code>\s\w\w\w\s</code> соответствует любому слову из трех букв, окруженному пробельными символами.
\S	Любой не пробельный символ.	Выражение <code>\s\S\S\S\s</code> соответствует любым трем непобельным символам, окруженным пробельными.
\d	Любая десятичная цифра	Выражение <code>c\d t</code> соответствует фрагментам: <code>c1t</code> , <code>c2t</code> , <code>c3t</code> и т.д.
\D	Любой символ, не являющийся десятичной цифрой	Выражение <code>c\D t</code> не соответствует фрагментам: <code>c1t</code> , <code>c2t</code> , <code>c3t</code> и т.д.

Кроме метасимволов, обозначающие классы символов, могут применяться уточняющие метасимволы:

Уточняющие символы	Описание
<code>^</code>	Фрагмент, совпадающий с регулярными выражениями, следует искать только в начале строки
<code>\$</code>	Фрагмент, совпадающий с регулярными выражениями, следует искать только в конце строки
<code>\A</code>	Фрагмент, совпадающий с регулярными выражениями, следует искать только в начале многострочной строки
<code>\Z</code>	Фрагмент, совпадающий с регулярными выражениями, следует искать только в конце многострочной строки
<code>\b</code>	Фрагмент, совпадающий с регулярными выражениями, начинается или заканчивается на границе слова, т.е. между символами, соответствующими метасимволам <code>\w</code> и <code>\W</code>
<code>\B</code>	Фрагмент, совпадающий с регулярными выражениями, не должен встречаться на границе слов

В регулярных выражениях часто используются повторители - метасимволы, которые располагаются непосредственно после обычного символа или группы символов и задают количество его повторений в выражении.

Повторители	Описание	Пример
<code>*</code>	Ноль или более повторений предыдущего элемента	Выражение <code>ca*t</code> соответствует фрагментам: <code>ct</code> , <code>cat</code> , <code>caat</code> , <code>caaat</code> и т.д.
<code>+</code>	Одно или более повторений предыдущего элемента	Выражение <code>ca+t</code> соответствует фрагментам: <code>cat</code> , <code>caat</code> , <code>caaat</code> и т.д.
<code>?</code>	Не более одного повторения предыдущего элемента	Выражение <code>ca?t</code> соответствует фрагментам: <code>ct</code> , <code>cat</code> .
<code>{n}</code>	Ровно <code>n</code> повторений предыдущего элемента	Выражение <code>ca{3}t</code> соответствует фрагменту: <code>caaat</code> . Выражение <code>(cat){2}</code> соответствует фрагменту: <code>catcat</code> .
<code>{n,}</code>	По крайней мере <code>n</code> повторений предыдущего элемента	Выражение <code>ca{3,}t</code> соответствует фрагментам: <code>caaaat</code> , <code>caaaaaaat</code> и т.д. Выражение <code>(cat){2,}</code> соответствует фрагментам: <code>catcat</code> , <code>catcatcat</code> и т.д.
<code>{n, m}</code>	От <code>n</code> до <code>m</code> повторений предыдущего элемента	Выражение <code>ca{2, 4}t</code> соответствует фрагментам: <code>caaat</code> , <code>caaaaat</code> .

Регулярное выражение записывается в виде строкового литерала, причем перед строкой необходимо ставить символ `@`, который говорит о том, что строку нужно будет рассматривать и в том случае, если она будет занимать несколько строчек на экране. Однако символ `@` можно не ставить, если в качестве шаблона используется шаблон без метасимволов.

**Замечание.** Если нужно найти какой-то символ, который является метасимволом, например, точку, можно это сделать защитив ее обратным слэшем. Т.е. просто точка означает любой одиночный символ, а `\.` означает просто точку.

Примеры регулярных выражений:

- слово `rus` - `@"rus"` или `"rus"`
- номер телефона в формате `xxx-xx-xx` - `@"\d\d\d-\d\d-\d\d"` или `@"\d{3}(-\d\d){2}"`
- номер автомобиля - `@"[A-Z]\d{3}[A-Z]{2}\d{2,3}RUS"`

**Задания.** Запишите регулярное выражение, соответствующее:

1. дате в формате дд.мм.гг или дд.мм.гггг
2. времени в формате чч.мм или чч:мм
3. целому числу (со знаком и без)
4. вещественному числу (со знаком и без, с дробной частью и без, с целой частью и без)

### Поиск в тексте по шаблону

Пространство имен библиотеки базовых классов `System.Text.RegularExpressions` содержит все объекты платформы .NET Framework, имеющие отношение к регулярным выражениям. Важнейшим классом, поддерживающим регулярные выражения, является класс `Regex`, который представляет неизменяемые откомпилированные регулярные выражения. Для описания регулярного выражения в классе определено несколько перегруженных конструкторов:

1. `Regex()` - создает пустое выражение;
2. `Regex(String)` - создает заданное выражение;
3. `Regex(String, RegexOptions)` - создает заданное выражение и задает параметры для его обработки с помощью элементов перечисления `RegexOptions` (например, различать или нет прописные и строчные буквы).

Поиск фрагментов строки, соответствующих заданному выражению, выполняется с помощью методов `IsMatch`, `Match`, `Matches` класса `Regex`.

Метод `IsMatch` возвращает `true`, если фрагмент, соответствующий выражению, в заданной строке найден, и `false` в противном случае. Например, попытаемся определить, встречается ли в заданном тексте слово `собака`:

```
static void Main()
{
    Regex r = new Regex("собака", RegexOptions.IgnoreCase);
    string text1 = "Кот в доме, собака в конуре.";
    string text2 = "Котик в доме, собачка в конуре.";
    Console.WriteLine(r.IsMatch(text1));
    Console.WriteLine(r.IsMatch(text2));
}
```

*Замечание.* `RegexOptions.IgnoreCase` - означает, что регулярное выражение применяется без учета регистра символов.

Можно использовать конструкцию выбора из нескольких элементов. Варианты выбора перечисляются через вертикальную черту. Например, попытаемся определить, встречается ли в заданном тексте слов `собака` или `кот`:

```
static void Main(string[] args)
{
    Regex r = new Regex("собака|кот", RegexOptions.IgnoreCase);
    string text1 = "Кот в доме, собака в конуре.";
    string text2 = "Котик в доме, собачка в конуре.";
    Console.WriteLine(r.IsMatch(text1));
    Console.WriteLine(r.IsMatch(text2));
}
```

Попытаемся определить, есть ли в заданных строках номера телефона в формате `xx-xx-xx` или `xxx-xx-xx`:

```
static void Main()
{
    Regex r = new Regex(@"\d{2,3}(-\d\d){2}");
    string text1 = "tel:123-45-67";
    string text2 = "tel:no";
    string text3 = "tel:12-34-56";
    Console.WriteLine(r.IsMatch(text1));
    Console.WriteLine(r.IsMatch(text2));
    Console.WriteLine(r.IsMatch(text3));
}
```

**Задание.** Измените программу так, чтобы можно было определить, содержится в тексте дата в формате дд.мм.гг.

Метод `Match` класса `Regex` не просто определяет, содержится ли текст, соответствующий шаблону, а возвращает объект класса `Match` - последовательность фрагментов текста, совпавших с шаблоном. Следующий пример позволяет найти все номера телефонов в указанном фрагменте текста:

```
static void Main()
{
    Regex r = new Regex(@"\d{2,3}(-\d\d){2}");
    string text = @"Контакты в Москве tel:123-45-67, 123-34-56; fax:123-56-45
                  Контакты в Саратове tel:12-34-56; fax:12-56-45";
    Match tel = r.Match(text);
}
```

```

while (tel.Success)
{
    Console.WriteLine(tel);
    tel = tel.NextMatch();
}
}

```

Следующий пример позволяет подсчитать сумму целых чисел, встречающихся в тексте:

```

static void Main()
{
    Regex r = new Regex(@"[-+]?[0-9]+");
    string text = @"5*10=50 -80/40=-2";
    Match teg = r.Match(text);
    int sum = 0;
    while (teg.Success)
    {
        Console.WriteLine(teg);
        sum += int.Parse(teg.ToString());
        teg = teg.NextMatch();
    }
    Console.WriteLine("sum=" + sum);
}

```

**Задание.** Измените программу так, чтобы на экран дополнительно выводилось количество найденных чисел.

Метод `Matches` класса `Regex` возвращает объект класса `MatchCollection` - коллекцию всех фрагментов заданной строки, совпавших с шаблоном. При этом метод `Matches` многократно запускает метод `Match`, каждый раз начиная поиск с того места, на котором закончился предыдущий поиск.

```

static void Main(string[] args)
{
    string text = @"5*10=50 -80/40=-2";
    Regex theReg = new Regex(@"[-+]?[0-9]+");
    MatchCollection theMatches = theReg.Matches(text);
    foreach (Match theMatch in theMatches)
    {
        Console.Write("{0} ", theMatch.ToString());
    }
    Console.WriteLine();
}
}

```

### *Редактирование текста*

Регулярные выражения могут эффективно использоваться для редактирования текста. Например, метод `Replace` класса `Regex` позволяет выполнять замену одного фрагмента текста другим или удаление фрагментов текста:

**Пример 1.** Изменение номеров телефонов:

```

static void Main(string[] args)
{
    string text = @"Контакты в Москве tel:123-45-67, 123-34-56; fax:123-56-45.
Контакты в Саратове tel:12-34-56; fax:11-56-45";
    Console.WriteLine("Старые данные\n"+text);
    string newText=Regex.Replace(text, "123-", "890-");
    Console.WriteLine("Новые данные\n" + newText);
}

```

**Задание.** Измените программу так, чтобы шестизначные номера заменялись на семизначные добавлением 0 после первых двух цифр. Например, номер 12-34-56 заменился бы на 120-34-56.

**Пример 2.** Удаление всех номеров телефонов из текста:

```

static void Main()
{
    string text = @"Контакты в Москве tel:123-45-67, 123-34-56; fax:123-56-45.
Контакты в Саратове tel:12-34-56; fax:12-56-45";
    Console.WriteLine("Старые данные\n"+text);
    string newText=Regex.Replace(text, @"\d{2,3}(-\d\d){2}", "");
    Console.WriteLine("Новые данные\n" + newText);
}
}

```

**Задание.** Измените программу так, чтобы из текста удалялись слова `tel` и `fax` (если после данных слов стоят двоеточия, то их тоже следует удалить).

**Пример 3.** Разбиение исходного текста на фрагменты:

```
static void Main()
{
    string text = @"Контакты в Москве tel:123-45-67, 123-34-56; fax:123-56-45.
                   Контакты в Саратове tel:12-34-56; fax:12-56-45";
    string []newText=Regex.Split(text,"[ ,.;;]+");
    foreach( string a in newText)
        Console.WriteLine(a);
}
```