

Практическая работа №14 Обработка строк. использование регулярных выражений в Джава приложениях

Цель: практической работы – понять особенности использования регулярных выражений в Java, научиться работать с строками и применять регулярные выражения для обработки строк в программах.

Теоретические сведения

Регулярные выражения и их использование в Java программах для обработки строк

Регулярные выражения – эта система обработки текста, основанная на специальной системе записи образцов для поиска. Образец (pattern), задающий правило поиска, по-русски также иногда называют «шаблоном», «маской». Сейчас регулярные выражения используются многими текстовыми редакторами и утилитами для поиска и изменения текста на основе выбранных правил. Язык программирования Java также поддерживает регулярные выражения для работы со строками.

Основными классами для работы с регулярные выражения являются класс `java.util.regex.Pattern` и класс `java.util.regex.Matcher`.

Для определения шаблона применяются специальные синтаксические конструкции.

Класс Pattern

Класс `java.util.regex.Pattern` применяется для определения регулярных выражений, для которого ищется соответствие в строке, файле или другом объекте представляющем собой некоторую последовательность символов. Этот класс используется для простой обработки строк. Для более сложной обработки строк используется класс `Matcher`, рассматриваемый ниже.

В классе `Pattern` объявлены следующие методы:

- `compile(String regex)` – возвращает `Pattern`, который соответствует `regex`;
- `matcher(CharSequence input)` – возвращает `Matcher`, с помощью которого можно находить соответствия в строке `input`;
- `matches(String regex, CharSequence input)` – проверяет на соответствие строки `input` шаблону `regex`;
- `pattern()` – возвращает строку, соответствующую шаблону;

- `split(CharSequence input)` – разбивает строку `input`, учитывая, что разделителем является шаблон;
- `split(CharSequence input, int limit)` – разбивает строку `input` на не более чем `limit` частей.

С помощью метода `matches()` класса `Pattern` можно проверять на соответствие шаблону целой строки, но если необходимо найти соответствия внутри строки, например, определять участки, которые соответствуют шаблону, то класс `Pattern` не может быть использован. Для таких операций необходимо использовать класс `Matcher`.

Класс `Matcher`

С помощью класса `java.util.regex.Matcher` можно получить больше информации каждом соответствии.

Начальное состояние объекта типа `Matcher` не определено. Попытка вызвать какой-либо метод класса для извлечения информации о найденном соответствии приведет к возникновению ошибки `IllegalStateException`. Для того чтобы начать работу с объектом `Matcher` нужно вызвать один из его методов:

- `matches()` – проверяет, соответствует ли вся строка шаблону;
- `lookingAt()` – пытается найти последовательность символов, начинающуюся с начала строки и соответствующую шаблону;
- `find()` или `find(int start)` – пытается найти последовательность символов, соответствующих шаблону, в любом месте строки.

Параметр `start` указывает на начальную позицию поиска.

Иногда необходимо сбросить состояние объекта класса `Matcher` в исходное, для этого применяется метод `reset()` или `reset(CharSequence input)`, который также устанавливает новую последовательность символов для поиска.

Для замены всех подпоследовательностей символов, удовлетворяющих шаблону, на заданную строку можно применить метод `replaceAll(String replacement)`.

Для того чтобы ограничить поиск границами входной последовательности применяется метод `region(int start, int end)`, а для получения значения этих границ – `regionEnd()` и `regionStart()`. С регионами связано несколько методов:

`useAnchoringBounds(boolean b)` – если установлен в `true`, то начало и конец региона соответствуют символам `^` и `$` соответственно;

`hasAnchoringBounds()` – проверяет закреплённость границ.

В регулярном выражении для более удобной обработки входной последовательности применяются группы, которые помогают выделить части найденной подпоследовательности. В шаблоне они обозначаются скобками «(» и «)». Номера групп начинаются с единицы. Нулевая группа совпадает со всей найденной подпоследовательностью. Далее приведены методы для извлечения информации о группах:

`end()` – возвращает индекс последнего символа подпоследовательности, удовлетворяющей шаблону;

`end(int group)` – возвращает индекс последнего символа указанной группы;

`group()` – возвращает всю подпоследовательность, удовлетворяющую шаблону;

`group(int group)` – возвращает конкретную группу;

`groupCount()` – возвращает количество групп;

`start()` – возвращает индекс первого символа подпоследовательности, удовлетворяющей шаблону;

`start(int group)` – возвращает индекс первого символа указанной группы;

`hitEnd()` – возвращает истину, если был достигнут конец входной последовательности.

Следующий пример показывает использование возможностей классов `Pattern` и `Matcher`, для поиска, разбора и разбиения строк.

Листинг 14.1 – Пример использования классов `Pattern` и `Matcher`,

```
import java.util.regex.*;

public class DemoRegular {

    public static void main(String[] args) {

        // проверка на соответствие строки шаблону

        Pattern p1 = Pattern.compile("a*y");

        Matcher m1 = p1.matcher( "aaay" );
```

```

boolean b = m1.matches();

System.out.println(b);

// поиск и выбор подстроки, заданной шаблоном
String regex = "(\\w+)@(\\w+\\.)(\\w+)(\\.\\w+)*" ;

String s = "адреса эл.почты: mymail@tut.by и
rom@bsu.by";

Pattern p2 = Pattern.compile(regex);
Matcher m2 = p2.matcher(s);

while (m2.find()) {
    System.out.println("e-mail: " + m2.group());
}

// разбиение строки на подстроки с применением
шаблона в качестве
// разделителя

Pattern p3 = Pattern.compile("\\d+\\s?");
String[] words = p3.split("java5tiger 77
java6mustang");

for (String word : words)
    System.out.println(word);
}
}}

```

В результате работы программы будет выведено:

```

true

e - mail : mymail @ tut. by

e-mail: rom@bsu.by

java

tiger

```

```
java
```

```
mustang
```

Следующий пример демонстрирует возможности использования групп, а также собственных и неполных квантификаторов.

Листинг 14.2 – Пример использования квантификаторов

```
import java.util.regex.*;

public class Groups {

    public static void main(String[] args) {
        String input = "abdcxyz";
        myMatches("([a-z]*) ([a-z]+)", input);
        myMatches("([a-z]?)([a-z]+)", input);
        myMatches("([a-z]+) ([a-z]*)", input);
    }
}
```

В результате работы программы будет выведено:

```
First group: abdcxy
```

```
Second group: z
```

```
First group: a
```

```
Second group: bdcxyz
```

```
First group: abdcxyz
```

```
Second group: nothing
```

В первом случае к первой группе (First group) относятся все возможные символы, но при этом остается минимальное количество символов для второй группы (Second group). Во втором случае для первой группы выбирается наименьшее количество символов, т. к. используется слабое совпадение. В третьем случае первой группе будет соответствовать вся строка, а для второй не остается ни одного символа, так как вторая группа использует слабое совпадение. В четвертом случае строка не соответствует регулярному выражению, т. к. для двух групп выбирается наименьшее количество символов.

В классе `Matcher` объявлены два полезных метода для замены найденных подпоследовательностей во входной строке.

`Matcher appendReplacement(StringBuffer sb, String replacement)` – метод читает символы из входной строки и добавляет их в `sb`. Чтение останавливается на `start()` – первой позиции предыдущего совпадения, после чего происходит добавление в `sb` строки `replacement`. При следующем вызове этого метода, производится добавление символов, начиная с символа с индексом `end()` предыдущего совпадения.

Класс `StringTokenizer`

Этот класс предназначен для разбора строки по лексемам (tokens). Строка которую необходимо разобрать передается в качестве параметра конструктору `StringTokenizer(String str)`. Определено еще два перегруженных конструктора, которым дополнительно можно передать строку-разделитель лексем `StringTokenizer(String str,String delim)` и признак возврата разделителя лексем `StringTokenizer(String str,String delim,Boolean returnDelims)`. Разделителем лексем по умолчанию служит пробел.

Листинг 14.3 – Пример класса `Test`

```
public class Test {  
  
    public Test() {  
    }  
  
    public static void main(String[] args) {  
        Test test = new Test();  
        String toParse = "word1;word2;word3;word4";  
        StringTokenizer st = new  
            StringTokenizer(toParse, ";");  
        while(st.hasMoreTokens()) {  
            System.out.println(st.nextToken());  
        }  
    }  
}
```

Вывод:

word1

word2

word3

word4

Задания на практическую работу №14

1. Необходимо реализовать консольное приложение, позволяющее манипулировать строкой, разбив ее на элементы путем использования регулярных выражений.
2. Написать регулярное выражение, определяющее является ли данная строка строкой "abcdefghijklmnpqrstuv18340" или нет.
 - а) пример правильных выражений: abcdefghijklmnpqrstuv18340
 - б) пример неправильных выражений: abcdefghijklmnoasdfsdpqrstuv18340.
3. Дан текст со списками цен. Извлечь из него цены в USD, RUB, EU.
 - а) пример правильных выражений: 25.98 USD.
 - б) пример неправильных выражений: 44 ERR, 0.004 EU.
4. Дан текст, необходимо проверить есть ли в тексте цифры, за которыми не стоит знак «+».
 - а) пример правильных выражений: $(1 + 8) - 9 / 4$
 - б) пример неправильных выражений: $6 / 5 - 2 * 9$
5. Написать регулярное выражение, определяющее является ли данная строка датой в формате dd/mm/yyyy. Начиная с 1900 года до 9999 года.
 - а) пример правильных выражений: 29/02/2000, 30/04/2003, 01/01/2003.
 - б) пример неправильных выражений: 29/02/2001, 30-04-2003, 1/1/1899.
6. Написать регулярное выражение, определяющее является ли данная строка допустимым (корректным) e-mail адресом согласно RFC под номером 2822.
 - а) пример правильных выражений: user@example.com, root@localhost
 - б) пример неправильных выражений: myhost@@@com.ru, @my.ru, Julia String.
7. Проверить, надежно ли составлен пароль. Пароль считается надежным, если он состоит из 8 или более символов. Где символом может быть цифр, английская буква, и знак подчеркивания. Пароль должен содержать хотя бы одну заглавную букву, одну маленькую букву и одну цифру.
 - а) пример правильных выражений: F032_Password, TrySpy1.

- b) пример неправильных выражений: smart_pass, A007.
8. Напишите метод filter, который принимает на вход массив (любого типа) и реализацию интерфейса Filter с методом apply(Object o), чтобы убрать из массива лишнее. Проверьте как он работает на строках или других объектах.
 9. Постройте частотный словарь¹ букв русского (или английского) алфавита. Здесь не имеет значения проблема выбора и анализа корпуса языка, достаточно будет взять текст небольшой длины).

¹https://ru.wikipedia.org/wiki/%D0%A7%D0%B0%D1%81%D1%82%D0%BE%D1%82%D0%BD%D1%8B%D0%B9_%D1%81%D0%BB%D0%BE%D0%B2%D0%B0%D1%80%D1%8C

