

Лабораторная работа №13 (3-й семестр)

Теоретические сведения

Измерять продуктивность программирования подсчетом строк кода — это так же, как оценивать постройку самолета по его весу.
— Bill Gates

Регулярные выражения (англ. regular expressions) — формальный язык поиска и осуществления манипуляций с подстроками в тексте, основанный на использовании метасимволов (символов-джокеров, англ. wildcard characters). По сути это строка-образец (англ. pattern, по-русски её часто называют «шаблоном», «маской»), состоящая из символов и метасимволов и задающая правило поиска.

Регулярные выражения используются некоторыми текстовыми редакторами и утилитами для поиска и подстановки текста. Например, при помощи регулярных выражений можно задать шаблоны, позволяющие:

- найти все последовательности символов «кот» в любом контексте, как то: «кот», «котлета», «терракотовый»;
- найти отдельно стоящее слово «кот» и заменить его на «кошка»;
- убрать из текста все предложения, в которых упоминается слово кот или кошка.

Регулярные выражения позволяют задавать и гораздо более сложные шаблоны поиска или замены.

Обычные символы (литералы) и специальные символы (метасимволы)

Большинство символов в регулярном выражении представляют сами себя за исключением специальных символов [] \ / ^ \$. | ? * + () { }, которые могут быть предварены символом \ (обратная косая черта) («экранированы», «защищены») для представления их самих в качестве символов текста. Можно экранировать целую последовательность символов, заключив её между \Q и \E.

Метасимвол . (точка) означает один любой символ, но в некоторых реализациях исключая символ новой строки.

Символьные классы (наборы символов)

Набор символов в квадратных скобках [] именуется символьным классом и позволяет указать интерпретатору регулярных выражений, что на данном месте в строке может стоять один из перечисленных символов. В частности, [абв] задаёт возможность появления в тексте одного из трёх указанных символов, а [1234567890] задаёт соответствие одной из цифр. Возможно указание диапазонов символов: например, [А-Яа-я] соответствует всем буквам русского алфавита, за исключением букв «Ё» и «ё».

Если требуется указать символы, которые не входят в указанный набор, то используют символ ^ внутри квадратных скобок, например [^0-9] означает любой символ, кроме цифр.

Некоторые символьные классы можно заменить специальными метасимволами:

Табл. 1 — Некоторые символьные классы

Символ	Эквивалент	Соответствие
\d	[0-9]	Цифровой символ
\D	[^0-9]	Нецифровой символ
\s	[\f\n\r\t\v]	Пробельный символ
\S	[^\f\n\r\t\v]	Непробельный символ
\w	[[:word:]]	Буквенный или цифровой символ или знак подчёркивания
\W	[^[:word:]]	Любой символ, кроме буквенного или цифрового символа или знака подчёркивания

Позиция внутри строки. Следующие символы позволяют позиционировать регулярное выражение относительно элементов текста: начала и конца строки, границ слова.

Табл. 2 — Метасимволы для позиционирования регулярного выражения в тексте

Представление	Позиция	Пример
<code>^</code>	Начало строки	<code>^a</code>
<code>\$</code>	Конец строки	<code>a\$</code>
<code>\b</code>	Граница слова	<code>a\b</code> <code>\ba</code>
<code>\B</code>	Не граница слова	<code>\Ba\B</code>
<code>\G</code>	Предыдущий успешный поиск	<code>\Ga</code>

Квантификация (поиск последовательностей). Квантификатор после символа, символьного класса или группы определяет, сколько раз предшествующее выражение может встречаться. Следует учитывать, что квантификатор может относиться более чем к одному символу в регулярном выражении, только если это символьный класс или группа.

Табл. 3 — Квантификаторы

Представление	Число повторений		Пример	Соответствие
<code>{n}</code>	Ровно n раз		<code>colou{3}r</code>	<code>colouuur</code>
<code>{m,n}</code>	От m до n включительно		<code>colou{2,4}r</code>	<code>colouur</code> , <code>colouuur</code> , <code>colouuuur</code>
<code>{m,}</code>	Не менее m		<code>colou{2,}r</code>	<code>colouur</code> , <code>colouuur</code> , <code>colouuuur</code> и т. д.
<code>{,n}</code>	Не более n		<code>colou{,3}r</code>	<code>color</code> , <code>colour</code> , <code>colouur</code> , <code>colouuur</code>
<code>?</code>	Ноль или одно	<code>{0,1}</code>	<code>colou?r</code>	<code>color</code> , <code>colour</code>
<code>*</code>	Ноль или более	<code>{0,}</code>	<code>colou*r</code>	<code>color</code> , <code>colour</code> , <code>colouur</code> и т. д.
<code>+</code>	Одно или более	<code>{1,}</code>	<code>colou+r</code>	<code>colour</code> , <code>colouur</code> и т. д. (но не <code>color</code>)

Часто используется последовательность `.*` для обозначения любого количества любых символов между двумя частями регулярного выражения.

Символьные классы в сочетании с квантификаторами позволяют устанавливать соответствия с реальными текстами. Например, столбцами цифр, телефонами, почтовыми адресами, элементами HTML-разметки и др.

Если символы `{ }` не образуют квантификатор, их специальное значение игнорируется.

Жадная и ленивая квантификация

В некоторых реализациях квантификаторам в регулярных выражениях соответствует максимально длинная строка из возможных (квантификаторы являются жадными, англ. *greedy*). Это может оказаться значительной проблемой. Например, часто ожидают, что выражение `(<.*>)` найдёт в тексте теги HTML. Однако, если в тексте есть более одного HTML-тега, то этому выражению соответствует целиком строка, содержащая множество тегов.

Выражение `(<.*>)` соответствует строке, содержащей несколько тегов HTML-разметки, целиком.

<p>Википедия — свободная энциклопедия, в которой <i>каждый</i> может изменить или дополнить любую статью.</p>

Чтобы выделить отдельные теги, можно применить ленивую версию этого выражения: (<.*?>) Ей соответствует не вся показанная выше строка, а отдельные теги (выделены цветом):

<p>Википедия — свободная энциклопедия, в которой <i>каждый</i> может изменить или дополнить любую статью.</p>

Использование ленивых квантификаторов может повлечь за собой обратную проблему, когда выражению соответствует слишком короткая, в частности, пустая строка.

Класс wxRegEx. Класс предоставляет инструменты для работы с регулярными выражениями. Класс находится в библиотеке «wx/regex.h». Основные члены класса приведены в табл.

Табл. - некоторые методы класса wxRegEx

Метод	Описание
wxRegEx(const wxString& expr, int flags = wxRE_DEFAULT)	Создает и компилирует регулярное выражение,
bool Compile(const wxString& pattern, int flags = wxRE_DEFAULT)	Компилирует регулярное выражение, возвращает истину в случае успеха
bool IsValid()	Возвращает истину, если регулярное выражение составлено без синтаксических ошибок
bool GetMatch(size_t* start, size_t* len, size_t index = 0) const	Возвращает индекс и длину вхождения выражения
size_t GetMatchCount() const	Возвращает длину массива вхождений
bool Matches(const wxChar* text, int flags = 0) const	Проверяет на наличие соответствий в введенном тексте
int ReplaceAll(wxString* text, const wxString& replacement) const	Заменяет в строке text все вхождения регулярного выражения на replacement

```
wxString text=_T("trololo ghjhf jhgfjhgf kghkhkgv a@a.com.ua");

wxRegEx re(_T(".*"));
if (re.IsValid())
{
    wxMessageBox(_T("Valid"));
    if (re.Matches(text))
    {
        wxString result = re.GetMatch(text, 0);
        wxMessageBox(result);
    }
}
```

Якоря		Образцы шаблонов	
^	Начало строки +	([A-Za-z0-9-]+)	Буквы, числа и знаки переноса
\A	Начало текста +	(\d{1,2}\V\d{1,2}\V\d{4})	Дата (напр., 21/3/2006)
\$	Конец строки +	([^\s]+(?:\.(jpg gif png))\.\w2)	Имя файла jpg, gif или png
\Z	Конец текста +	(^[1-9]{1}\$ ^[1-4]{1}[0-9]{1}\$ ^[50]\$)	Любое число от 1 до 50 включительно
\b	Граница слова +	(#?([A-Fa-f0-9]){3}([A-Fa-f0-9]){3})?)	Шестнадцатиричный код цвета
\B	Не граница слова +	((?=[*\d])(?=[*a-z])(?=[*A-Z]).{8,15})	От 8 до 15 символов с минимум одной цифрой, одной заглавной и одной строчной буквой (полезно для паролей).
\<	Начало слова	(\w+@[a-zA-Z_]+?\.[a-zA-Z]{2,6})	Адрес email
\>	Конец слова	(\</?[^\>]+\>)	HTML теги
Символьные классы		Примечание Эти шаблоны предназначены для ознакомительных целей и основательно не проверялись. Используйте их с осторожностью и предварительно тестируйте.	
\c	Управляющий символ		
\s	Пробел	Кванторы	
\S	Не пробел		
\d	Цифра	Диапазоны	
\D	Не цифра		
\w	Слово	Примечание Диапазоны включают граничные значения.	
\W	Не слово		
\xhh	Шестнадцатиричный символ hh	Модификаторы шаблонов	
\Oxxx	Восьмиричный символ xxx		
Символьные классы POSIX		Мета-символы (экранируются)	
[upper:]	Буквы в верхнем регистре		
[lower:]	Буквы в нижнем регистре	Примечание Эта таблица доступна на www.exlab.net Англоязычный оригинал на AddedBytes.com	
[alpha:]	Все буквы		
[alnum:]	Буквы и цифры	Подстановка строк	
[digit:]	Цифры		
[xdigit:]	Шестнадцатиричные цифры	Утверждения	
[punct:]	Пунктуация		
[blank:]	Пробел и табуляция	Примечание Отмеченное + работает в большинстве языков программирования.	
[space:]	Пустые символы		
[cntrl:]	Управляющие символы	Примечание	
[graph:]	Печатные символы		
[print:]	Печатные символы и пробелы	Примечание	
[word:]	Буквы, цифры и подчеркивание		