



**Нижегородский государственный университет
им. Н.И.Лобачевского**

Институт Информационных Технологий, Математики и Механики

Регулярные выражения в C++ (Regex)

Содержание

- Регулярные выражения
- Синтаксис
- Примеры



Регулярные выражения

- **Регулярные выражения** (англ. *regular expressions*) — формальный язык поиска и осуществления манипуляций с подстроками в тексте, основанный на использовании метасимволов (символов-джокеров, англ. *wildcard characters*). Для поиска используется строка-образец (англ. *pattern*, по-русски её часто называют «шаблоном», «маской»), состоящая из символов и метасимволов и задающая правило поиска. Для манипуляций с текстом дополнительно задаётся строка замены, которая также может содержать в себе специальные символы.



Регулярные выражения

■ Функции регулярных выражений:

- ❑ проверка наличия искомого образца в заданном тексте;
- ❑ определение подстроки текста, которая сопоставляется образцу;
- ❑ определение групп символов, соответствующих отдельным частям образца.

■ Пример использования регулярных выражений:

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



Регулярные выражения в C++

- Стандартная библиотека в C++11.
В Visual Studio 2008 пространство имен `std::tr1::regex`, в более поздних `std::regex`.
- Поддерживаемые синтаксисы:
 - ☐ basic
 - ☐ extended
 - ☐ ECMAScript (по умолчанию в VS)
 - ☐ awk
 - ☐ grep
 - ☐ egrep
 - ☐ и т.д.



Регулярные выражения в C++

■ Основные методы класса `regex`:

- ❑ **`basic_regex (const CharT* s, flag_type f = std::regex_constants::ECMAScript);`** - конструктор, `CharT* s` – шаблон, `flag_type f` – тип синтаксиса

■ Функции с регулярными выражениями:

- ❑ **`bool regex_match(string str, regex rx);`** - функция сравнивает строку `str` с шаблоном `rx`. Если строка соответствует шаблону, то возвращает `true`, иначе `false`
- ❑ **`bool regex_search(string str, smatch match, regex rx);`** - функция ищет в строке `str` все соответствующие шаблону `rx` вхождения и записывает их в `match`. Возвращает `true` если найдено хотябы одно соответствие
- ❑ **`string regex_replace(string str, regex rx, string replace, ...);`** - функция заменяет все вхождения шаблона `rx` в `str` на `replace` и возвращает полученную строку.
- ❑ Возможны иные прототипы у каждой из функций



Синтаксис ECMAScript

- В грамматике ECMAScript специальное значение имеют следующие символы:
 - `^ $ \ . * + ? () [] { } |`
- Остальные символы считаются обычными.



Синтаксис ECMAScript

- «\» - экранирующий символ, после него символ интерпретируется по другому. Специальный становится обычным, а обычный интерпретируется как специальный.
- «^» - соответствие началу ввода. Например: «^A» соответствует «A» в «An E», но не в «an A».
- «\$» - соответствует концу ввода. Пример шаблона: «A\$».
- Повторение предыдущего символа:
 - «*» - 0 и более раз;
 - «+» - 1 и более раз;
 - «?» - 0 или 1 раз;
 - «{n}» - ровно n вхождений;
 - «{n,m}» - вхождений от n до m раз;
 - Пример «ba+» соответствует «ba», «baa», «baaa» и тд, но не «b».
- «.» - один любой символ, кроме переноса строки.
- «[...]» - соответствие из набора символов Пример: [0-9] или [0123456789] – соответствие цифре. [^0-9] – не цифра.
- «|» - или.
- «(...)» - захватывающие скобки



Синтаксис ECMAScript

- `\b` - соответствует границе слова.
- `\B` - соответствует несловообразующей границе.
Несловообразующая граница соответствует позиции, в которой предыдущий и следующий символы являются символами одного типа.
- `\d` – соответствует цифровому символу.
- `\D` – соответствует нецифровому символу.
- `\n`, `\r` – соответствуют переводу строки и возврату каретки.
- `\s` – одиночный символ пустого пространства. Пробел, `tab` и т.д.
- `\S` – одиночный символ непустого пространства.
- `\t` – горизонтальный `tab`.
- `\v` – вертикальный `tab`.
- `\w` – любой цифробуквенный символ, включая нижнее подчеркивание. Эквивалент `[A-Za-z0-9_]`.
- `\W` – любой не цифробуквенный символ.
- И т.д.



Пример

- `(2[0-3]|[0-1]\d):[0-5]\d` – время в формате «часы:минуты»
- `^((8|\+7)(\[-])?)(\(?\d{3}\)?(\[-])?)[\d\[-]{7,10}$` – номер телефона в любом доступном формате:
 - ❑ `(8|\+7)` – номер начинается с 8 или +7
 - ❑ `[\[-]?` – разделение в виде «-» или « », которое может отсутствовать.
 - ❑ `((8|\+7)(\[-])?` – начальный префикс может отсутствовать.
`(495)123456`
 - ❑ `\(?\d{3}\)?` – наличие трех цифр, в скобках и без.
 - ❑ `[\d\[-]{7,10}` – наличие от 7 до 10 цифр, с разделениями

