正则表达式快速参考手册

胡志飞 <WisdomFusion#gmail.com>

2012年6月2日

目 录

1	简介	1
2	基本语法	1
3	高级语法	6
4	举些栗子	8
5	正则表达式"流派"	8
6	应用场景	9
	6.1 正则表达式工具箱	ç
	6.2 应用案例	C

1 简介 Introduction

正则表达式,(Regular Expression, 在代码中常简写为 regex、regexp 或 RE),计算机科学的一个概念。正则表达式使用字符来描述、匹配一系列符合某个句法规则的字符串。在很多文本编辑器里,正则表达式通常被用来检索、替换那些符合某个模式的文本。许多程序设计语言都支持利用正则表达式进行字符串操作。例如,在 Perl¹中就内建了一个功能强大的正则表达式引擎。正则表达式这个概念最初是由 Unix 中的工具软件(例如 sed²和 grep³)普及开的。

需要注意的是,用什么工具,用什么编辑语言,正则表达式的语法有些差别,特性的支持也参差不齐,称之为正则表达式"流派"(第5部分详述),所以要单独参考工具和编程语言本身的文档才行。

2 基本语法 Basic Syntax

特性	语法	描述	举个栗子
字符	除 [\^\$. ?*+() 以 外的任意字符	除了[\^\$. ?*+()以外的任意字符, {和}也是文字文本,除了下面说到的成对出现的量词语法,如 {n}和 {m,n}等。	a 匹配 about 中的 a
	字符转义	\t,\?,*,\+,\.,\ ,\{,\},\,\[,\],\(,\)	\+ 匹配 + ; \?\- 匹配 ?-
	\n,\r 和 \t	Windows 文件格式换行符是\r\n, UNIX 文件格式换行符是\n, \t	

To be continued...

¹Perl被称为"实用报表提取语言"(Practical Extraction and Report Language),正则表达式特性的推动者,文本处理非常方便。

²sed是一种 UNIX/Linux 平台下的轻量级流编辑器,日常一般用于处理文本文件。

³grep,global search regular expression and print out the line,是一种强大的文本搜索工具,它能使用正则表达式搜索文本,并把匹配的行打印出来。

特性 语法		描述	举个栗子
表符 (\xoB)		Ctrl+A 到 Ctrl+Z,Ctrl+a 到 Ctrl+z	
		依次为警报(\xo7)、Esc字符(\x1B)、进纸符(\xoC)和垂直制表符(\xoB)	
		文字文本范围,被包含在 \Q 和 \E 之间的文字,都被视为普通文字,如 [\^\$. ?*+(){} 也不再用转义了,这个最早是由 Perl 引入正则表达式的。	\Q+-*/\E 匹配的就是 +-*/
基本特性	基本特性 · (点) 匹配除换行符之外的任意字符,有些正则表达式"流派"还支 是否匹配换行符的开关。		. 匹配 about 中的任意一个字符
	I	管道,或的关系,匹配 的左侧或右侧的字符串	abc def xyz 匹配 abc 或 def 或 xyz
字符类	[]	匹配字符类中列举的任意一个字符	[abc] 配 a 或 b 或 c ; [.!?] 配 . 或!或?
	[\^\]]	在字符类中,要匹配 ^-]\这几字符,得使用\转义	[\^\]] 匹配 ^ 或]
	[^]	排除型字符类, ^ (脱字符, caret) 紧跟 [之后, 可以把字符类中列举的字符排除匹配范围, 也就是所这个字符类将匹配任意一个不在列出字符范围内的字符	

To be continued...

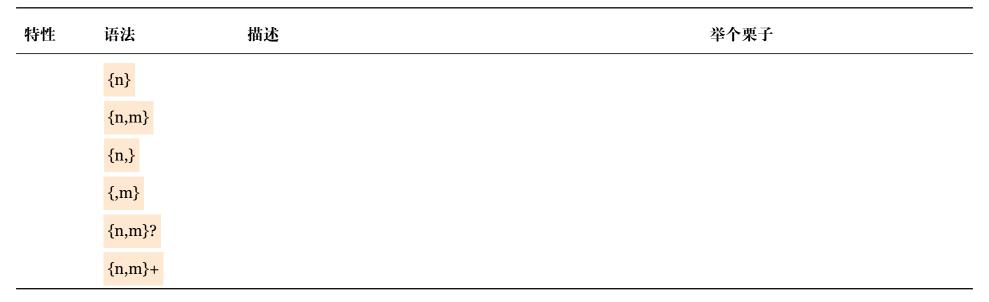
特性	语法	描述	举个栗子
	\d, \w, \s	\d 匹配数字,与 [o-9] 等价; \w 匹配任意一个字母或数字或下划线或汉字; \s 匹配任意一个空白符	[\d\s] 匹配一个数字或空白符
	\D, \W, \S	是 \d, \w 和 \s 的反义字符类。 \D 匹配任意非数字的字符; \W 匹配任意不是字母、数字、下划线、汉字的字符; \S 匹配任意不是空白符的字符	\D 匹配任意非数字的字符
	[\b]	在字符类中, [\b] 为 Backspace 退格键字符	
POSIX	[:alnum:]		
	[:alpha:]		
	[:ascii:]		
	[:blank:]		
	[:cntrl:]		
	[:digit:]		
	[:graph:]		
	[:lower:]		

To be continued...

特性	语法	描述	举个栗子	
	[:print:]			
	[:punct:]			
	[:space:]			
	[:upper:]			
	[:word:]			
	[:xdigit:]			
锚点	^			
	\$ \A			
	\A			
	\G			
	\Z			
	\G \Z \b \B			
	\<			

特性	语法	描述	举个栗子
	\>		
	\`		
	\` \& \'		
	\'		
量词	??		
	??		
	? +		
	*		
	* *? *+ + +?		
	*+		
	+		
	++		

To be continued...



3 高级语法 Advanced Syntax

特性	语法	描述	举个栗子
Unicode			

特性	语法	描述	举个栗子

 分组与反向引用
 (regex)

 (1到)9
 (10到)99

 (g{1}到)g{99}
 (g{-1}, \g{-2}, etc.

 (?<name>regex)
 (k<name>, \g{name}

 (g{name}
 (g{name}

To be continued...

特性	语法	描述	举个栗子
高级分组	(?#comment)		
	(?=Regex)		
	(?!Regex)		
	(?<=regex)		
	(? regex)</td <td></td> <td></td>		

4 举些栗子 Regex Examples

一些栗子

5 正则表达式"流派"Regex Flavors

6 应用场景 Application Scenarios

6.1 正则表达式工具箱 Regex Toolbox

总有一款适合你, Windows 下的记事本太鸡肋, Word 处理方式主要是"通配符"而不是正则表达式。

RegexBuddy
平台的实现和增强。

JGsoft 开发的一个强大的正则表达式测试工具, UltraEdit, Notepad++

grep Vim

PowerGREP GNU Emacs

RegexBuddy 的兄弟软件, 同是 JGsoft 开发, 是 grep 在 Windows sed & awk

6.2 应用案例 Application Cases

Dreamweaver 表格处理

VBA 中使用正则表达式

```
Sub IndentParaWithRegEx()

PowerPoint VBA 挑量给指定字符开头段落加动画
Dim oSld As Slide
Dim oShp As Shape
Dim i As Integer

Full 和关变量
Dim regx As Object, oMatch As Object

Sub IndentParaWithRegEx()

PowerPoint VBA 挑量给指定字符开头段落加动画
Dim oSld As Slide
Dim oShp As Shape
Dim i As Integer

Full 和关变量
Dim regx As Object, oMatch As Object

Sub Teg 查找的正则,参考 http://msdn.microsoft.com/en-us/library/ms974570.aspx
strPattern = "~开头字符串"

Set regx = CreateObject("vbscript.regexp")
With regx

Global = True
```

10 6 应用场景

```
.IgnoreCase = True
15
       .Pattern = strPattern
16
   End With
18
   For Each oSld In ActivePresentation.Slides
       For Each oShp In oSld.Shapes
20
           If oShp.HasTextFrame Then
21
               If oShp.TextFrame2.HasText Then
22
                   With oShp.TextFrame2.TextRange
23
                       For i = 1 To .Paragraphs.Count
24
                           With .Paragraphs(i)
                                ' 可能会出现多个匹配项的
26
                               If (regx.Test(.Text) = True) Then
27
                                    .ParagraphFormat.FirstLineIndent = 0
28
                               End If
29
                           End With
30
                       Next i 'para
31
                   End With
32
               End If 'has text
33
           End If 'has textframe
34
       Next oShp
   Next oSld
   End Sub
```

InDesign GREP

使用 Perl 正则表达式处理文件 神的编辑器之正则