

正则表达式快速参考手册

胡志飞

<WisdomFusion[at]gmail[dot]com>

2012 年之旧文重拾于 2016 年 2 月 26 日
version 0.1.0

目 录

1	简介	1	6.2.1	Perl 正则表达式王国	13
2	基本语法	2	6.2.2	Python	13
3	高级语法	8	6.2.3	PHP PCRE	13
4	举些栗子	11	6.2.4	sed & awk	13
5	正则表达式“流派”	12	6.2.5	grep	13
6	应用场景	12	6.2.6	Swift	13
6.1	正则表达式工具箱	12	6.2.7	Java	13
6.2	应用案例	13	6.2.8	JavaScript	13
			6.2.9	Adobe Dreamweaver 表格处理	13
			6.2.10	VBA 中使用正则表达式	13
			6.2.11	GREP for Adobe InDesign	15
			6.2.12	神的编辑器 GNU Emacs 之正则神器	15
			7	后记	15

1 简介 Introduction

文字处理无处不在无时不有，日常工作和学习大多数任务都和文字息息相关，编辑们写文章、整理资料，开发人员编码、处理用户提交的数据或请求接口数据，等等，这些都是以字符和字符串相关的任务，既然如此，掌握一个快速文字处理的方法就变得很有必要。

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 **Perl**^①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 **sed**^②和 **grep**^③）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。本文档旨在给大家一个通用的、概括的正则表达式宏观印象，辅以实例和应用案例，同时针对个别常用但又不易理解的特性，给大家作详细说明和总结，抛砖引玉。

期望本文档能给大家一个快速的参考，快速的掌握正则表达式这个棒棒哒效率工具，让大家平时工作学习中更加得心应手！☺

说明

我对排版及专业出版知之甚微，只是在平时笔记和文档时，哪怕是自己写的太乱的话也不乐意翻看，印象笔记里的东东又过于零碎，故把旧文完善并整排^a。然而，专业领域知识因涉猎过多而不精，难保周全和准确，但只要在自己知识圈内，我会劲力完成尽可能规范和可靠的文档呈现给大家，并不断完善更新，请发现问题后联系指正，共同提高。

^a本文档 2012 年编写，现整拾，并加以完善。目前觉得 Markdown, **org-mode**适用于文本文档的快速编写，而 Adobe InDesign 和 LaTeX 适合更专业的文档和书籍的图文混排及设计。

^①**Perl**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。

^②**sed**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。

^③**grep**, 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。

2 基本语法 Basic Syntax

语法部分结合了自己的理解和对正则表达式应用的一些心得，分类有不当之处，请指正。因为是总结性的参考文档，所以这里使用“大表哥”形式展示，列出语法的同时，关键语法举了几个栗子加强理解。

特性	语法	描述	举个栗子
字符	除 <code>[^\\$. ?*\+0]</code> 以外的任意字符	除了 <code>[^\\$. ?*\+0]</code> 以外的任意字符， <code>{}</code> 和 <code>}</code> 也是文字文本，除了下面说到的成对出现的量词语法，如 <code>{n}</code> 和 <code>{m,n}</code> 等。	<code>a</code> 匹配 <code>about</code> 中的 <code>a</code>
字符转义	<code>\t</code> , <code>\?</code> , <code>*</code> , <code>\+</code> , <code>\.</code> , <code>\ </code> , <code>\{</code> , <code>\}</code> , <code>\\</code> , <code>\[</code> , <code>\]</code> , <code>\(</code> , <code>\)</code> <code>\n</code> , <code>\r</code> 和 <code>\t</code> <code>\cA</code> 到 <code>\cZ</code> , <code>\ca</code> 到 <code>\cz</code> <code>\a</code> , <code>\e</code> , <code>\f</code> , <code>\v</code> <code>\Q ... \E</code>	<code>\t</code> , <code>\?</code> , <code>*</code> , <code>\+</code> , <code>\.</code> , <code>\ </code> , <code>\{</code> , <code>\}</code> , <code>\\</code> , <code>\[</code> , <code>\]</code> , <code>\(</code> , <code>\)</code> Windows 文件格式换行符是 <code>\r\n</code> , UNIX 文件格式换行符是 <code>\n</code> , <code>\t</code> 匹配水平制表符 <code>Ctrl</code> + <code>A</code> 到 <code>Ctrl</code> + <code>Z</code> , 与 ASCII 字符 <code>\x01</code> 到 <code>\x1A</code> 等价 依次为警报 (<code>\x07</code>), <code>Esc</code> 字符 (<code>\x1B</code>)、进纸符 (<code>\x0C</code>) 和垂直制表符 (<code>\x0B</code>) 文字文本范围，被包含在 <code>\Q</code> 和 <code>\E</code> 之间的文字，都被视为普通文字，如 <code>[^\\$. ?*\+0]</code> 也不再转义了，这个最早是由 Perl 引入正则表达式的。	<code>\+</code> 匹配 <code>+</code> ; <code>\?\-</code> 匹配 <code>?-</code> <code>\Q+*\E</code> 匹配的就是 <code>+*/</code>
基本特性	<code>.</code> (点)	匹配除换行符之外的任意字符，有些正则表达式“流派”还支持点是否匹配换行符的开关。	<code>.</code> 匹配 <code>about</code> 中的任意一个字符

To be continued...

(续表)

特性	语法	描述	举个栗子
		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] 匹配 a 或 b 或 c [aeiou] 匹配任何一个英文元音字母 [.!?] 匹配 . 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	[^a-d] 匹配除了 a,b,c,d 之外的任意一个字符
	\d, \w, \s	\d 匹配数字，与 [0-9] 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:]	匹配所有大小写字母及数字	等价于 [0-9a-zA-Z]

To be continued...

(续表)

特性	语法	描述	举个栗子
	<code>[:alpha:]</code>	匹配所有大小写字母	等价于 <code>[a-zA-Z]</code>
	<code>[:ascii:]</code>	匹配所有 ASCII 字符, 查看完整 ASCII 字符列表	等价于 <code>[\x01-\x7F]</code>
	<code>[:blank:]</code>	匹配半角空格和制表符	等价于 <code>[\t]</code>
	<code>[:cntrl:]</code>	匹配所有 ASCII 0 到 31 之间的控制符	等价于 <code>[\x01-\x1F]</code>
	<code>[:digit:]</code>	匹配所有数字	等价于 <code>[0-9]</code>
	<code>[:graph:]</code>	匹配所有可打印的字符	
	<code>[:lower:]</code>	匹配所有小写字母	等价于 <code>[a-z]</code>
	<code>[:print:]</code>	匹配所有可打印字符和空格	
	<code>[:punct:]</code>	匹配所有标点符号	
	<code>[:space:]</code>	空白字符	等价于 <code>[\t\n\r\f\v]</code>
	<code>[:upper:]</code>	匹配所有大写字母	等价于 <code>[A-Z]</code>
	<code>[:word:]</code>	字母、数字和下划线	等价于 <code>[a-zA-Z0-9_]</code>
	<code>[:xdigit:]</code>	匹配所有十六进制字符	等价于 <code>[0-9a-fA-F]</code>

To be continued...

(续表)

特性	语法	描述	举个栗子
锚点	<code>^</code>	匹配字符串开始位置或行首位置	单行模式下 <code>^.</code> 在 <code>foo\nbar</code> 中匹配 <code>f</code> ; 在多行模式下, 同时还匹配换行后的 <code>b</code>
	<code>\$</code>	匹配字符串结尾位置或行尾位置	<code>.\$</code> 在 <code>foo\nbar</code> 中匹配 <code>r</code> ; 在多行模式下, 同时还匹配换行符前的 <code>o</code>
	<code>\A</code>	字符串开头位置 (类似 <code>^</code> , 但不受处理多行选项的影响)	<code>\Ae</code> 在 <code>example</code> 这个字符串中匹配开头的 <code>e</code>
	<code>\Z</code>	字符串结尾位置或行尾位置 (不受处理多行选项的影响)	<code>e\Z</code> 在 <code>example</code> 这个字符串中匹配结尾的 <code>e</code>
	<code>\b</code>	单词分界位置, 单词开头或结尾	<code>.\b</code> 在字符串 <code>abc</code> 中匹配 <code>c</code>
	<code>\B</code>	匹配不是单词开头或结尾的位置	<code>\B.\B</code> 在字符串 <code>abc</code> 中匹配 <code>b</code>
	<code>\<</code>	单词开头	
	<code>\></code>	单词结尾	

To be continued...

(续表)

特性	语法	描述	举个栗子
量词	?	前导字符重复零次或一次，贪婪的 ^① ：当正则表达式中包含能接受重复的限定符时，通常的行为是（在使整个表达式能得到匹配的前提下）匹配尽可能多的字符。	abc? 匹配 abc 或 ab，如果可能，优先匹配前者
	??	前导字符重复零次或一次，非贪婪 ^② ：当正则表达式中包含能接受重复的限定符时，通常的行为是（在使整个表达式能得到匹配的前提下）匹配尽可能少的字符。与贪婪相反。	abc?? 匹配 ab 或 abc
	*	前导字符重复零次或更多次，贪婪的	
	*?	前导字符重复零次或更多次，非贪婪	
	+	前导字符重复一次或更多次，贪婪的	
	+	前导字符重复一次或更多次，非贪婪	
	{n}	前导字符重复 n 次	

To be continued...

^①贪婪模式：当正则表达式中包含能接受重复的限定符时，通常的行为是（在使整个表达式能得到匹配的前提下）匹配尽可能多的字符。考虑这个表达式：`a.*b`，它将会匹配最长的以 `a` 开始，以 `b` 结束的字符串。如果用它来搜索 `aabab` 的话，它会匹配整个字符串 `aabab`。这被称为贪婪匹配。

^②相反，非贪婪，即匹配尽可能少的字符，只要在量词后面加上一个问号`?`。如 `a.*?b` 匹配最短的，以 `a` 开始，以 `b` 结束的字符串。如果把它应用于 `aabab` 的话，它会匹配 `aab`（第一到第三个字符）和 `ab`（第四到第五个字符）。

(续表)

特性	语法	描述	举个栗子
	{n,m}	前导字符重复 n 到 m 次，其中 $n \geq 0$ ， $m \geq n$	
	{n,}	前导字符重复 n 次或更多次，其中 $n \geq 0$	
	{,m}	前导字符最多重复 m 次，基中 $m \geq 0$	
分组与 反向引用	(regex)	匹配 regex，并捕获文本到自动命名的组里	(abc){3} 匹配 abcabcabc
	(?:regex)	匹配 regex，不捕获匹配的文本，也不给此分组分配组号	(?:abc){3} 匹配 abcabcabc，无分组
	\1 到 \9	反向引用，用于重复搜索前面某个分组匹配的文本。例如，\1 代表分组 1 匹配的文本。有些与此正则表达式流派支持多于 9 的分组	(abc def)=\1 匹配 abc=abc 或 def=def，而不是 abc=def 或 def=abc
	\10 到 \99	反向引用，分组 10 到 99	
	\g{1} 到 \g{99}	Perl 语法中，反向引用语法优化 ^①	避免出现歧义，同时用负数分组还能倒序引用。
	\g{-1}, \g{-2}, etc.	倒数第 1 个分组，倒数第 2 个分组， ...	
	(?<name>regex)	命名分组	命名分组的最大好处是反向引用时不用再怕弄错分组了。

To be continued...

^①如果想实现类似 (.)\1 的效果，若紧跟的字符和分组号相同，如 (.)\111，这样正则表达式引擎就不知所措，如果用 \g{1}11 这种语法就不会有歧义，最重要的是这种语法或以用负数分组号倒序选用分组。

(续表)

特性	语法	描述	举个栗子
	<code>\k<name></code>	反向引用命令分组，Perl 中也可以使用 <code>\g{name}</code> 。	
	<code>\`</code>	正则表达式匹配部分之前的字符串，Perl 语言中也作 <code>\${^PREMATCH}</code>	
	<code>\&</code>	正则表达式匹配的部分，Perl 语言中也作 <code>\${^MATCH}</code>	
	<code>\'</code>	正则表达式匹配部分之后的字符串，Perl 语言中也作 <code>\${^POSTMATCH}</code>	

3 高级语法 Advanced Syntax

之所以本文档中称这些正则为“高级语法”，一是有些不常用，二是有些语法不太好理解，故有此一说。

特性	语法	描述	举个栗子
模式 修饰符	<code>(?i)</code>	打开忽略大小写模式，之后的模式不分大小写。模式修饰符有 <code>i</code> , <code>s</code> , <code>m</code> , <code>x</code> 四种，分别是忽略大小写 (<code>IgnoreCase</code>) 、单行模式 (<code>Singleline</code>) 、多行模式 (<code>Multiline</code>) 和注释模式	<code>(?i)te(?-i)st</code> 匹配 <code>TEst</code> ，而不匹配 <code>TEST</code>

To be continued...

(续表)

特性	语法	描述	举个栗子
	(?-i)	关闭忽略大小写模式，之后的模式不分大小写	
	(?s)	打开单行模式，之后的模式不支持多行	默认情况下 . 是不匹配换行的，打开该模式后，待匹配的字符将为视为“一行”。
	(?-s)	关闭单行模式，之后的模式支持多行	
	(?m)	打开多行模式，之后的模式支持多行	
	(?-m)	^ 和 \$ 匹配行首和行尾	
	(?x)	打开宽松和注释模式	打开该模式后，可以在正则表达式中插和空白和换行，使正则表达式可读性增强。
	(?-x)	关闭宽松和注释模式	
	(?i-sm)	打开 i 和 m 模式，关闭 s 模式	以上几种模式可以组合使用
	(?i-sm:regex)	在 (?i-sm:regex) 子模式内打开 i 和 m 模式，关闭 s 模式	
注释	(?#comment)	注释	
零宽断言	(?=Regex)	肯定顺序环视 (Negative Lookahead) 子表达式 能够匹配 右侧 文本，含以下三种些类正则被称“环视 (Lookaround, Lookahead 和 Lookbehind 统称为 Lookaround)”，也称为“零宽断言”(Zero-Length Assertions)	\b\w+(?=ing\b) , 匹配以 ing 结尾的单词的前面部分 (除了 ing 以外的部分)

To be continued...

(续表)

特性	语法	描述	举个栗子
	(?!Regex)	否定顺序环视 (Positive Lookahead) 子表达式 <u>不能</u> 匹配 <u>右侧</u> 文本	<code>(?<=\bre)\w+\b</code> 会匹配以 re 开头的单词的后半部分 (除了 re 以外的部分) 又, <code>(?<=\s)\d+(?=\s)</code> 匹配以空白符间隔的数字 (再次强调, 不包括这些空白符)
	(?<=regex)	肯定逆序环视 (Positive Lookbehind) 子表达式 <u>能</u> 够匹配 <u>左侧</u> 文本	<code>\d{3}(?!\d)</code> 匹配三位数字, 而且这三位数字的后面不能是数字
	(?<!\d)	否定逆序环视 (Negative Lookbehind) 子表达式 <u>不能</u> 匹配 <u>左侧</u> 文本	<code>(?<![a-z])\d{7}</code> 匹配前面不是小写字母的七位数字 又, <code>(?<=<(\w+)>).(?=<\/\1>)</code> 匹配不包含属性的简单 HTML 标签内里的内容
固化分组	(?>regex)	贪婪子表达式, 也称“ 固化分组 ”, 使用它可以加快匹配失败的速度, 如 Subject 这个字符串, 现用 <code>^\w:</code> 对其进行匹配, 正则表达式引擎发现 Subject 不匹配, 就会试图匹配 Subjec , 一直尝试到 S , 发现都不匹配才得出无法匹配的结论。如果使用固化分组 <code>^(?>\w+):</code> , 它会直接试图使用 <code>\w+</code> 去匹配 Subjec 字符串, 而不会一一回溯, 发现 <code>\w+</code> 后面没有 <code>:</code> , 立即报告失败。	如果字符串中没有第二个 x 的时候, <code>x(?>\w+)x</code> 要比 <code>x\w+x</code> 高效得多

4 举些例子吧！Regex Examples

Email 地址

`^[a-zA-Z0-9_+-.]+@[a-zA-Z0-9]+\.[a-zA-Z0-9-.]++$` 匹配形如 `WisdomFusion@gmail.com` 的邮箱地址。

日期

`^\d{4}\-(0?[1-9]|1[012])\-(0?[1-9]|12[0-9]|3[01])$` 匹配形如 `yyyy-mm-dd` 格式 的日期。

非零负整数

`^\-[1-9][0-9]*$` 匹配如 `-2`, `-1024`, ... 之类的非零负整数。

匹配浮点数

`[-+]?([0-9]*\.[0-9]+|[0-9]+)` 匹配如 `3.1415926` 的浮点数 (带小数位的)。

去除重复行

查找 `^(.*) (\r?\n\1)+$`, 替的为 `\1`。

匹配用户名

`^[a-z0-9_-]{3,16}$`

网址

`^(https?:\/\/)?([\da-z\.-]+)\.([a-z\.] {2,6})([\/\w \.-]*)*\/?$`

IPv4 地址

`^((?:25[0-5]|2[0-4][0-9]|01?[0-9][0-9]?)\.){3}(?:25[0-5]|2[0-4][0-9]|01?[0-9][0-9]?)$`

HTML 标记

`<([a-z]+)([^\<]+)*(?:>(.*)<\/\1>|\s+\/>)`

UBB 代码清理

把\[/?(?:font|size)([^\]]+)?\[替换为空。

汉字

[^u4E00-u9FA5] 或 [一-]

5 正则表达式“流派” Regex Flavors

流派

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。

JGsoft RegexBuddy

JGsoft 开发的一个强大的正则表达式测试工具，这款是正则测试界最强大的工具了，没有之一，要墙裂向大家推介的哦！☺

JGsoft PowerGREP

PowerGREP RegexBuddy 的兄弟软件，同是 JGsoft 开发，是 grep 在 Windows 平台的实现和增强。

Debuggex

<https://www.debuggex.com/>

grep

grep

UltraEdit, Notepad++

UltraEdit, Notepad++

Vim

编辑器之神 Vim

GNU Emacs

神的编辑器 Emacs

sed & awk

sed and awk

6.2 应用案例 Application Cases

6.2.1 Perl 正则表达式王国

6.2.2 Python

6.2.3 PHP PCRE

6.2.4 sed & awk

6.2.5 grep

6.2.6 Swift

6.2.7 Java

6.2.8 JavaScript

6.2.9 Adobe Dreamweaver 表格处理

6.2.10 VBA 中使用正则表达式

```
1 Sub IndentParaWithRegEx()  
2 ' PowerPoint VBA 批量给指定字符开头段落加动画  
3 Dim oSld As Slide  
4 Dim oShp As Shape  
5 Dim i As Integer  
6 ' 正则相关变量  
7 Dim regx As Object, oMatch As Object  
8  
9 ' 这里写查找的正则, 参考 http://msdn.microsoft.com/en-us/library/ms974570.aspx  
10 strPattern = "^开头字符串"  
11  
12 Set regx = CreateObject("vbscript.regexp")
```

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


6.2.11 GREP for Adobe InDesign

6.2.12 神的编辑器 GNU Emacs 之正则神器

7 后记

后记