

1 正则表达式

```
#导入re包
import re

#function to check whether the pattern p can be found somewhere inside the string
s.
#re.search(p, s)

# $这个符号用于找到以ed结尾的词
[w for w in wordlist if re.search('ed$', w)]
#['abaissed', 'abandoned', 'abased', 'abashed', 'abatished', 'abed', 'aborted',
... ]

#通配符.指代单一字符
[w for w in wordlist if re.search('^.j...t..$', w)]
#['abjectly', 'adjuster', 'dejected', 'dejectly', 'injector', 'majestic', ...]
```

the ? symbol specifies that the previous character is optional. Thus «ê-?mail\$» will match both email and e-mail.

```
#方括号中间的字母表示寻找的字符可以是任意其中之一
[w for w in wordlist if re.search('^[ghi][mno][jlk][def]$', w)]
#['gold', 'golf', 'hold', 'hole']

# +号表示可以重复前面的字符或方括号中的字母
>>> [w for w in chat_words if re.search('^m+i+n+e+$', w)]
# ['miiiiiiiiiiiiinnnnnnnnnnneeeeeeee', 'miiiiinnnnnnnnneeeeeeee', 'mine', '
mmmmmmmmiiiiiiiiinnnnnnnnneeeeeeee']
>>> [w for w in chat_words if re.search('^ha+$', w)]
# ['a', 'aaaaaaaaaaaaaaaa', 'aaahhhh', 'ah', 'ahah', 'ahahah', 'ahh', 'ahhahahaha',
'ahhh', 'ahhhh', 'ahhhhhh', 'ahhhhhhhhhhhhh', 'h', 'ha', 'haaa', 'hah', '
haha', 'hahaaa', 'hahah', 'hahaha', 'hahahaa', 'hahahah', 'hahahaha', ...]
```

It should be clear that + simply means "one or more instances of the preceding item", which could be an individual character like m, a set like [fed] or a range like [d-f]. Now let's replace + with *, which means "zero or more instances of the preceding item". The regular expression «m*i*n*e*\$» will match everything that we found using «m+i+n+e+\$», but also words where some of the letters don't appear at all, e.g. me, min, and mmmmm. Note that the + and * symbols are sometimes referred to as Kleene closures, or simply closures.

字符`^`表示以之后的字符开头的，此外还有另外一种用法，就是在方括号的开头，比如《`[^aeiouAEIOU]`》就是找不是元音的字母

```
>>> wsj = sorted(set(nltk.corpus.treebank.words()))
>>> [w for w in wsj if re.search('^[0-9]+\.[0-9]+$', w)]
#['0.0085', '0.05', '0.1', '0.16', '0.2', '0.25', '0.28', '0.3', '0.4',
   '0.5', '0.50', '0.54', '0.56', '0.60', '0.7', '0.82', '0.84', '0.9', '0.95',
   '0.99', '1.01', '1.1', '1.125', '1.14', '1.1650', '1.17', '1.18', '1.19',
   '1.2', ...]
>>> [w for w in wsj if re.search('^[A-Z]+\$', w)]
#['C$', 'US$']
>>> [w for w in wsj if re.search('^[0-9]{4}$', w)]
#['1614', '1637', '1787', '1901', '1903', '1917', '1925', '1929', '1933', ...]
>>> [w for w in wsj if re.search('^[0-9]+-[a-z]{3,5}$', w)]
#['10-day', '10-lap', '10-year', '100-share', '12-point', '12-year', ...]
>>> [w for w in wsj if re.search('^[a-z]{5}-[a-z]{2,3}-[a-z]{,6}$', w)]
#['black-and-white', 'bread-and-butter', 'father-in-law', 'machine-gun-toting',
   'savings-and-loan']
>>> [w for w in wsj if re.search('(ed|ing)$', w)]
#['62%-owned', 'Absorbed', 'According', 'Adopting', 'Advanced', 'Advancing', ...]
```

Table 3.3:

Basic Regular Expression Meta-Characters, Including Wildcards, Ranges and Closures

Operator	Behavior
<code>.</code>	Wildcard, matches any character
<code>^abc</code>	Matches some pattern <i>abc</i> at the start of a string
<code>abc\$</code>	Matches some pattern <i>abc</i> at the end of a string
<code>[abc]</code>	Matches one of a set of characters
<code>[A-Z0-9]</code>	Matches one of a range of characters
<code>ed ing s</code>	Matches one of the specified strings (disjunction)
<code>*</code>	Zero or more of previous item, e.g. <i>a*</i> , <i>[a-z]*</i> (also known as <i>Kleene Closure</i>)
<code>+</code>	One or more of previous item, e.g. <i>a+</i> , <i>[a-z]+</i>
<code>?</code>	Zero or one of the previous item (i.e. optional), e.g. <i>a?</i> , <i>[a-z]?</i>
<code>{n}</code>	Exactly <i>n</i> repeats where <i>n</i> is a non-negative integer
<code>{n,}</code>	At least <i>n</i> repeats
<code>{,n}</code>	No more than <i>n</i> repeats
<code>{m,n}</code>	At least <i>m</i> and no more than <i>n</i> repeats
<code>a(b c)+</code>	Parentheses that indicate the scope of the operators