

Table 13.6 Regular Expression Object Methods

Syntax	Description
<code>rx.findall(s start, end)</code>	Returns all nonoverlapping matches of the regex in string <i>s</i> (or in the <i>start: end</i> slice of <i>s</i>). If the regex has captures, each match is returned as a tuple of captures.
<code>rx.finditer(s start, end)</code>	Returns a match object for each nonoverlapping match in string <i>s</i> (or in the <i>start: end</i> slice of <i>s</i>)
<code>rx.flags</code>	The flags that were set when the regex was compiled
<code>rx.groupindex</code>	A dictionary whose keys are capture group names and whose values are group numbers; empty if no names are used
<code>rx.match(s, start, end)</code>	Returns a match object if the regex matches at the start of string <i>s</i> (or at the start of the <i>start: end</i> slice of <i>s</i>); otherwise, returns <code>None</code>
<code>rx.pattern</code>	The string from which the regex was compiled
<code>rx.search(s, start, end)</code>	Returns a match object if the regex matches anywhere in string <i>s</i> (or in the <i>start: end</i> slice of <i>s</i>); otherwise, returns <code>None</code>
<code>rx.split(s, m)</code>	Returns the list of strings that results from splitting string <i>s</i> on every occurrence of the regex doing up to <i>m</i> splits (or as many as possible if no <i>m</i> is given). If the regex has captures, these are included in the list between the parts they split.
<code>rx.sub(x, s, m)</code>	Returns a copy of string <i>s</i> with every (or up to <i>m</i> if given) match replaced with <i>x</i> —this can be a string or a function; see text
<code>rx.subn(x, s m)</code>	The same as <code>re.sub()</code> except that it returns a 2-tuple of the resultant string and the number of substitutions that were made

One common task is to take an HTML text and output just the plain text that it contains. Naturally we could do this using one of Python’s parsers, but a simple tool can be created using regexes. There are three tasks that need to be done: delete any tags, replace entities with the characters they represent, and insert blank lines to separate paragraphs. Here is a function (taken from the `html2text.py` program) that does the job:

```
def html2text(html_text):
    def char_from_entity(match):
        code = html.entities.name2codepoint.get(match.group(1), 0xFFFD)
        return chr(code)
```