

Python读写XML文档(xml方式)

发布时间：2017年4月6日 18:05 作者：杨仕航
分类标签：Python 阅读(18654) 评论(1)

读写xml文档可以说是python基础。前段时间一直用xml.dom的方式读写xml文档。
但用起来比较繁琐，尤其是读写文本的时候特别麻烦。寻思换一种新的方式读写xml。

这个lxml和xml.etree.ElementTree两个的操作方式看起来差不多。xml.etree.ElementTree也是读写xml文档的一种方式。

lxml要更好一些，使用更简洁。解析xml的时候，自动处理各种编码问题。而且它天生支持 XPath 1.0。XSLT 1.0。定制元素类。

不过，lxml不是Python自带的标准库。需要自己安装。如下方式安装：

```
1. pip install lxml
```

下面看看如何使用lxml读写xml文档。

1、读取xml文档

1) 文档解析

lxml可以解析xml的字符串，使用etree.fromstring方法。如下所示：

```
1. #coding:utf-8
2. from lxml import etree
3.
4. xml_text = '<xml><head></head><body></body></xml>'
5. xml = etree.fromstring(xml_text)
```

lxml可以直接读取xml文件。

为了演示方便，找个xml文档作为例子。文件名为test.xml：

```
1. <?xml version="1.0" encoding="utf-8"?>
2. <root version="1.2" tag="test">
3.   <head>
4.     <title>test xml document</title>
5.   </head>
6.   <body>
7.     <items id="1">
8.       <source>aa</source>
9.     <target>AA</target>
10.   </items>
11.   <items id="2">
12.     <source>bb</source>
13.   <target>BB</target>
14.   </items>
15.   <items id="3">
16.     <source>cc</source>
17.   <target id="3t"><Cbpt id="3t1"/><cc</target>
18.   </items>
19. </body>
20. </root>
```

lxml读取xml文件的代码如下所示：

```
1. #coding:utf-8
2. from lxml import etree
3.
4. xml = etree.parse('test.xml') #读取test.xml文件
```

2) 获取属性

根节点root中有两个属性，我们可以通过如下方法获取根节点和其属性：

```
1. #coding:utf-8
2. from lxml import etree
3.
4. xml = etree.parse('test.xml') #读取test.xml文件
5. root = xml.getroot() #获取根节点
6.
7. #获取属性
8. print(root.items()) #获取全部属性和属性值
9. print(root.keys()) #获取全部属性
10. print(root.get('version', '')) #获取具体某个属性
```

得到如下结果：

```
1. [(('version', '1.0'), ('tag', 'test'))]
2. [('version', 'tag')]
3. 1.2
```

3) 获取节点

假如我们不知道root节点下有什么节点，可以通过循环遍历。

```
1. for node in root.getchildren():
2.     print(node.tag) #输出节点的标签名
```

得到如下结果：

```
1. head
2. body
```

若xml文档比较大，还可以使用itergetchildren方法。该方法得到一个生成器。
这里，你可以用dir(root)可以查看节点对象有什么方法。可以获取兄弟节点、父节点等方法。

若我们知道文档结构，还可以使用XPath方法获取节点。

XPath是一种表达式，可以快速查找XML文档中的信息。可以查看w3school教程。

不要被吓到，该语法很多简单。边查阅边使用，很快可以上手。w3school中也举了一些实例：

实例

在下面的表格中，我们已列出了一些路径表达式以及表达式的结果：

路径表达式	结果
bookstore	选取 bookstore 元素的所有子节点。
/bookstore	选取根元素 bookstore。 注释：路径必须起始于正斜杠 (/)，因此路径始终代表到根元素的绝对路径！
bookstore/book	选取属于 bookstore 的子元素的所有 book 元素。
//book	选取所有 book 子元素，而不管它们在文档中的位置。
bookstore//book	选取属于 bookstore 元素的所有后代的所有 book 元素，而不管它们位于 bookstore 之下的什么位置。
//@lang	选取名为 lang 的所有属性。

杨仕航的博客 <http://yshblog.com>

那么，我要获取全部items节点。使用XPath方法得到一个元素列表。代码如下：

```
1. root.xpath('//items')
```

可得到3个items元素。无论items元素的位置在哪。

我们也可以通过root和items之间的相对路径，严格得到我们想要位置的元素：

```
1. root.xpath('body/items')
```

一样可以得到3个元素，不过该xpath语法要严格很多。

假如head节点下也有items元素，该语法就不会获取到该元素。

若指定的XPath语法获取不到任何元素，将返回空列表。

4) 获取文本

有些元素中有文本，这个可以通过text属性获取。

```
1. #获取source元素中的文本
2. for node in root.xpath('//source'):
3.     print(node.text)
```

另外，还有一种情况。节点和文本混合的情况。

如test.xml中id为3t的target元素。target元素中有两段文本，以及文本中间还有个bpt元素。

在xml.minidom的方法，文本也是一种节点。这个问题好解决。

而在lxml中，文本不是节点。这种情况需要通过itertext方法或tail属性解决。

先尝试获取该元素的text属性：

```
1. #获取id属性为3t的target元素，这里后面的[0]
2. target = root.xpath('//target[@id="3t"]')[0]
3.
4. #输出该元素的text属性值
5. print(target.text)
```

将得到“CC”，后面的节点和“cc”获取不到。则text属性是获取到该节点下的第1段文本。

若该节点先是一个节点，再是文字：

```
1. <target id="3t"><bpt id="3t1"/><Cccc</target>
```

text属性将为None。

我们可以用itertext方法获取全部文本：

```
1. ''.join(target.itertext())
```

将得到“Cccc”，itertext方法得到一个生成器。该生成器是该节点下的全部文本生成器列表。

那假如我还需要获取其中的子节点。确保xml的结构的话，就需要使用tail属性。

tail属性是获取节点后的文本。我们可以先用text属性获取第1个文本，其他文本都通过子节点的tail属性获取。例如，我需要获取上面target元素的全部文本，若碰到子节点，则获取其id属性值一起拼成一个字符串。

```
1. texts = []
2.
3. #获取第1段文本
4. if target.text:
5.     texts.append(target.text)
6.
7. #遍历子节点
8. for sub in target.iterchildren():
9.     texts.append('bpt: ' + sub.get('id', ''))
10.    texts.append(sub.tail)
11.
12. #合并结果
13. print(''.join(texts))
```

将得到“CC-3t1-cc”。

2、写入xml文档

说完读取xml文档，进入写入xml文档环节。

现我们可以以模板test.xml文件，用lxml创建新的xml文档。结构和test.xml文件一致。

1) 创建文档（节点）

对于lxml来说，任意节点都可以保存成一个xml文档。

我们只需要给该节点加入属性、内容、子节点等即可。

那么创建节点方法如下：

```
1. #coding:utf-8
2. from lxml import etree
3.
4. #创建标签为root的节点
5. root = etree.Element('root')
```

在创建节点的同时，也可以给该节点加入命名空间：

```
1. root = etree.Element('root', namespace='xmlns':'http://www.w3.org/1999/xhtml')
```

在上面的test.xml中，还有两组属性。可用set方法添加属性：

```
1. root.set('version', '1.2')
2. root.set('tag', 'test')
```

当然，也可以在创建节点的时候，就写入属性：

```
1. attribs = {'version':'1.2', 'tag':'test'}
2. root = etree.Element('root', attrib=attribs)
```

2) 添加子节点

添加根节点之后，根节点下有两个子节点：head和body。

添加子节点有两种方法。先看方法1：

```
1. head = etree.Element('head')
2. root.append(head)
```

该方法是创建节点，再用append方法追加到root节点中。

还有一种方法，直接创建子节点：

```
1. head = etree.SubElement(root, 'head')
```

推荐使用第2种方法，比较快捷。

若需要写属性值，除了set方法。etree.SubElement方法也可以像etree.Element方法一样直接写入属性。

```
1. head = etree.SubElement(root, 'head', attrib=({'id':'head_id'})
```

3) 添加文本

test.xml文档中，有几个地方需要添加文本。先给head添加title属性，并加入文本：

```
1. title = etree.SubElement(head, 'title')
2. title.text = 'test xml document'
```

直接给text赋值即可。

其他文本写入我就忽略不写了，比较简单。

比较复杂的情况是上面提到的节点和文本混合的情况。这里同样给tail属性赋值即可。

```
1. body = etree.SubElement(root, 'body')
2. items = etree.SubElement(body, 'items', attrib=({'id':'3'})
3. target = etree.SubElement(items, 'target', attrib=({'id':'3t'})
4.
5. #写入第1段文本
6. target.text = 'CC'
7.
8. #写入CC后面的节点
9. bpt = etree.SubElement(target, 'bpt', attrib=({'id':'3t1'})
10.
11. #写入第2段文本，即bpt元素后面的文本
12. bpt.tail = 'cc'
```

通过，该方法可以有条件的写入文本。

4) 保存文档

文档写好之后，就保存文档。保存文档这里有两种方法。

一种为通过etree.tostring方法得到xml的文本，再手动写入。这个方法过于麻烦，就不讲了，也不推荐。

常规方法是通过etree的tree对象保存文件。代码如下：

```
1. #节点转为tree对象
2. tree = etree.ElementTree(root)
3. tree.write('test.xml', pretty_print=True, xml_declaration=True, encoding='utf-8')
```

各个参数含义如下：

- 1) 第1个参数是xml的完整路径(包括文件名)；
- 2) pretty_print参数是否美化代码；
- 3) xml_declaration参数是否写入xml声明，就是我们看到xml文档第1行文字；
- 4) encoding参数很明显是保存的编码。

lxml方式对比xml.minidom方式，有没觉得方便很多。

(原创博文，转载请注明出处 杨仕航的博客！本文链接：<http://yshblog.com/blog/151>)

若对你有帮助，不妨扫一扫右侧的二维码打赏支持我 ^_^

分享到：

上一篇：Excel的日期并不特殊

下一篇：我的网站搭建(第51天) 解除评论发送邮件限制

评论列表

我要评论

智慧如你，不想发表一下意见吗？

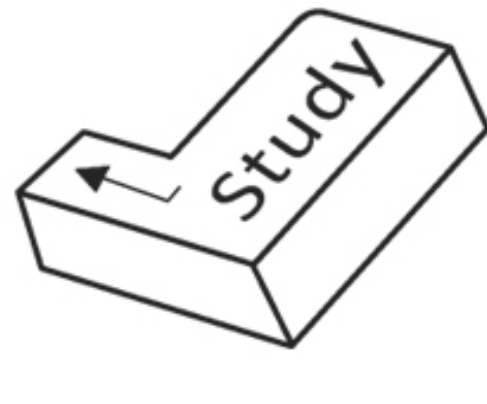
新的评论

为了防止垃圾评论，需要登录，才能评论哦~

😊🐼🍏🌻🏆🚀📺🚩

提交 清空

我的官网



再敲一行代码

Just one more line of code

随机推荐

我的网站搭建(第24天) 阅读计数优化
Django一些有用的admin设置
我的网站搭建(第46天) 在线头像
html最后一个元素不样式
Access建表规范信息玩(6)：最小信息
vba正则表达式入门
Python读写docx文件
Python用win32api操作注册表
我的网站搭建(第45天) 上传头像
我的网站搭建(第32天) OAuth功能整合

猜你喜欢

Python用win32api操作注册表
Python如何开发桌面软件
Python按位运算
我对Python装饰器的理解
Python写个猜数字游戏一玩
Python字符串处理方法总结
Python的json模块基本用法
Python求解一元二次方程
Python轻松实现排列和组合
Python2和3兼容写法

关于本站

- 1、基于Django+Bootstrap开发
- 2、主要发表本人的技术原创博客
- 3、本站于 2015-12-01 开始建站

建议反馈

- 1、可在相应的博文底下评论
- 2、发邮件到2872402050@qq.com

友情链接

再敲一行代码
QQ群：701914136