

xpath与json

有同学说，我正则用的不好，处理HTML文档很累，有没有其他的方法？

有！那就是XPath，我们可以先将 HTML文件 转换成 XML文档，然后用 XPath 查找 HTML 节点或元素。

什么是XML

- XML 指可扩展标记语言（EXtensible Markup Language）
- XML 是一种标记语言，很类似 HTML
- XML 的设计宗旨是传输数据，而非显示数据
- XML 的标签需要我们自行定义。
- XML 被设计为具有自我描述性。
- XML 是 W3C 的推荐标准

W3School官方文档：<http://www.w3school.com.cn/xml/index.asp>

XML 和 HTML 的区别

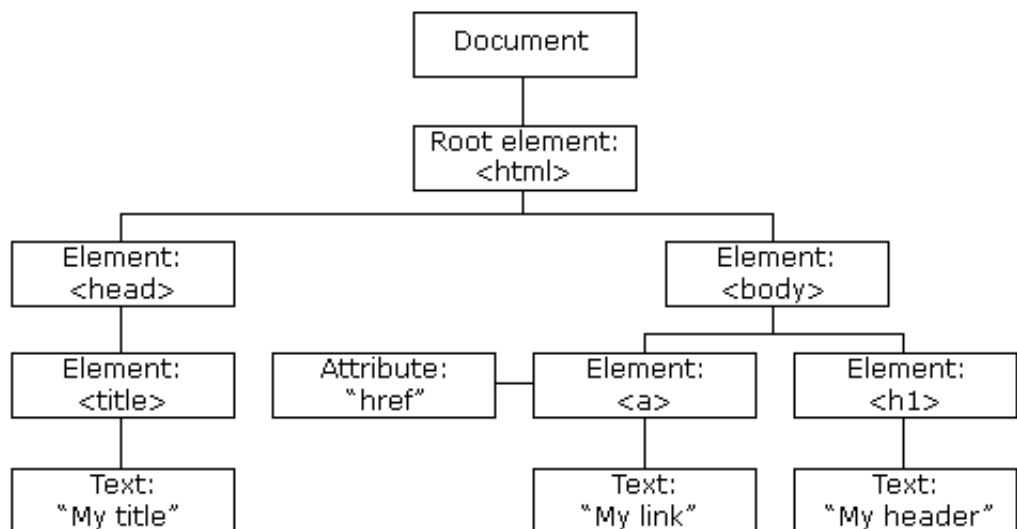
数据格式	描述	设计目标
XML	Extensible Markup Language (可扩展标记语言)	被设计为传输和存储数据，其焦点是数据的内容。
HTML	HyperText Markup Language (超文本标记语言)	显示数据以及如何更好显示数据。
HTML DOM	Document Object Model for HTML (文档对象模型)	通过 HTML DOM，可以访问所有的 HTML 元素，连同它们所包含的文本和属性。可以对其中的内容进行修改和删除，同时也可以创建新的元素。

```
1  <?xml version="1.0" encoding="utf-8"?>
2
3  <bookstore>
4
5      <book category="cooking">
6          <title lang="en">Everyday Italian</title>
7          <author>Giada De Laurentiis</author>
8          <year>2005</year>
9          <price>30.00</price>
10     </book>
11
12     <book category="children">
13         <title lang="en">Harry Potter</title>
14         <author>J K. Rowling</author>
15         <year>2005</year>
```

```
16     <price>29.99</price>
17 </book>
18
19 <book category="web">
20     <title lang="en">XQuery Kick Start</title>
21     <author>James McGovern</author>
22     <author>Per Bothner</author>
23     <author>Kurt Cagle</author>
24     <author>James Linn</author>
25     <author>Vaidyanathan Nagarajan</author>
26     <year>2003</year>
27     <price>49.99</price>
28 </book>
29
30 <book category="web" cover="paperback">
31     <title lang="en">Learning XML</title>
32     <author>Erik T. Ray</author>
33     <year>2003</year>
34     <price>39.95</price>
35 </book>
36
37 </bookstore>
```

HTML DOM 模型示例

HTML DOM 定义了访问和操作 HTML 文档的标准方法，以树结构方式表达 HTML 文档。



什么是XPath?

XPath (XML Path Language) 是一门在 XML 文档中查找信息的语言，可用在 XML 文档中对元素和属性进行遍历。

W3School官方文档: <http://www.w3school.com.cn/xpath/index.asp>

Xpath常用语法

表达式	说明	示例
/	从根节点开始选取	/html/body/div[3]/div/div[1]/div[2]/div[1]/div/h2
//	选择所有指定的结点	//div
.	当前结点	./div
..	父节点	../div
@	选取属性	//a/@href
text()	获取标签内的字符	//li/text()
node/elem[i]	选取node结点下的第i个elem元素 注意：下标从1开始。	ul/li[2]
node/elem[last()]	选取node结点下的最后一个elem元素	ul/li[last()]
//elem[@attr='xxx']	选取所有属性attr为xxx的elem元素	//div[@class='content']

初步使用

我们利用它来解析 HTML 代码，简单示例：

```
1 # 使用 lxml 的 etree 库
2 from lxml import etree
3
4 text = '''
5 <div>
6     <ul>
7         <li class="item-0"><a href="link1.html">first
8 item</a></li>
9         <li class="item-1"><a href="link2.html">second
10 item</a></li>
11         <li class="item-inactive"><a
12 href="link3.html">third item</a></li>
13         <li class="item-1"><a href="link4.html">fourth
14 item</a></li>
15         <li class="item-0"><a href="link5.html">fifth
16 item</a> # 注意，此处缺少一个 </li> 闭合标签
17     </ul>
18 </div>
19 '''
20
21 #利用etree.HTML，将字符串解析为HTML文档
22 html = etree.HTML(text)
23
24 # 按字符串序列化HTML文档
25 result = etree.tostring(html)
26
27 print(result)
```

输出结果：

```
1 <html><body>
2 <div>
3     <ul>
4         <li class="item-0"><a href="link1.html">first
item</a></li>
5         <li class="item-1"><a href="link2.html">second
item</a></li>
6         <li class="item-inactive"><a
href="link3.html">third item</a></li>
7         <li class="item-1"><a href="link4.html">fourth
item</a></li>
8         <li class="item-0"><a href="link5.html">fifth
item</a></li>
9     </ul>
10 </div>
11 </body></html>
```

lxml 可以自动修正 html 代码，例子里不仅补全了 li 标签，还添加了 body，html 标签。

文件读取：

除了直接读取字符串，lxml还支持从文件里读取内容。我们新建一个 hello.html 文件：

```
1 <!-- hello.html -->
2
3 <div>
4     <ul>
5         <li class="item-0"><a href="link1.html">first
item</a></li>
6         <li class="item-1"><a href="link2.html">second
item</a></li>
7         <li class="item-inactive"><a href="link3.html">
<span class="bold">third item</span></a></li>
8         <li class="item-1"><a href="link4.html">fourth
item</a></li>
9         <li class="item-0"><a href="link5.html">fifth
item</a></li>
10     </ul>
11 </div>
```

再利用 `etree.parse()` 方法来读取文件。

```
1 from lxml import etree
2
3 # 读取外部文件 hello.html
4 html = etree.parse('./hello.html')
5 result = etree.tostring(html, pretty_print=True)
6
7 print(result)
```

输出结果与之前相同：

```

1 <html><body>
2 <div>
3     <ul>
4         <li class="item-0"><a href="link1.html">first
item</a></li>
5         <li class="item-1"><a href="link2.html">second
item</a></li>
6         <li class="item-inactive"><a
href="link3.html">third item</a></li>
7         <li class="item-1"><a href="link4.html">fourth
item</a></li>
8         <li class="item-0"><a href="link5.html">fifth
item</a></li>
9     </ul>
10 </div>
11 </body></html>

```

XPath实例测试

<1> 获取所有的标签

```

1 from lxml import etree
2
3 html = etree.parse('hello.html')
4 print type(html) # 显示etree.parse() 返回类型
5
6 result = html.xpath('//li')
7
8 print(result) # 打印<li>标签的元素集合
9 print(len(result))
10 print(type(result))
11 print(type(result[0]))

```

<2> 继续获取标签的所有 class属性


```
1 from lxml import etree
2
3 html = etree.parse('hello.html')
4 result = html.xpath('//li/@class')
5
6 print(result)
```

运行结果

```
1 ['item-0', 'item-1', 'item-inactive', 'item-1', 'item-0']
```

<3> 获取 标签下的所有标签

```
1 from lxml import etree
2
3 html = etree.parse('hello.html')
4
5 #result = html.xpath('//li/span')
6 #注意这么写是不对的:
7 #因为 / 是用来获取子元素的, 而 <span> 并不是 <li> 的子元素, 所以, 要用双斜杠
8
9 result = html.xpath('//li//span')
10
11 print(result)
```

运行结果

```
1 [<Element span at 0x10d698e18>]
```

<4> 获取最后一个的<a>的href

```
1 from lxml import etree
2
3 html = etree.parse('hello.html')
4
5 result = html.xpath('//li[last()]/a/@href')
6 # [last()] 可以找到最后一个元素
7
8 print(result)
```

运行结果

```
1 ['link5.html']
```

<5> 获取倒数第二个元素的内容

```
1 from lxml import etree
2
3 html = etree.parse('hello.html')
4 result = html.xpath('//li[last()-1]/a')
5
6 # text 方法可以获取元素内容
7 print(result[0].text)
```

运行结果

```
1 fourth item
```

JSON格式数据与json模块

JSON(JavaScript Object Notation) 是一种轻量级的数据交换格式，它使得人们很容易的进行阅读和编写。同时也方便了机器进行解析和生成。适用于进行数据交互的场景，比如网站前台与后台之间的数据交互。

JSON和XML的比较可谓不相上下。

JSON

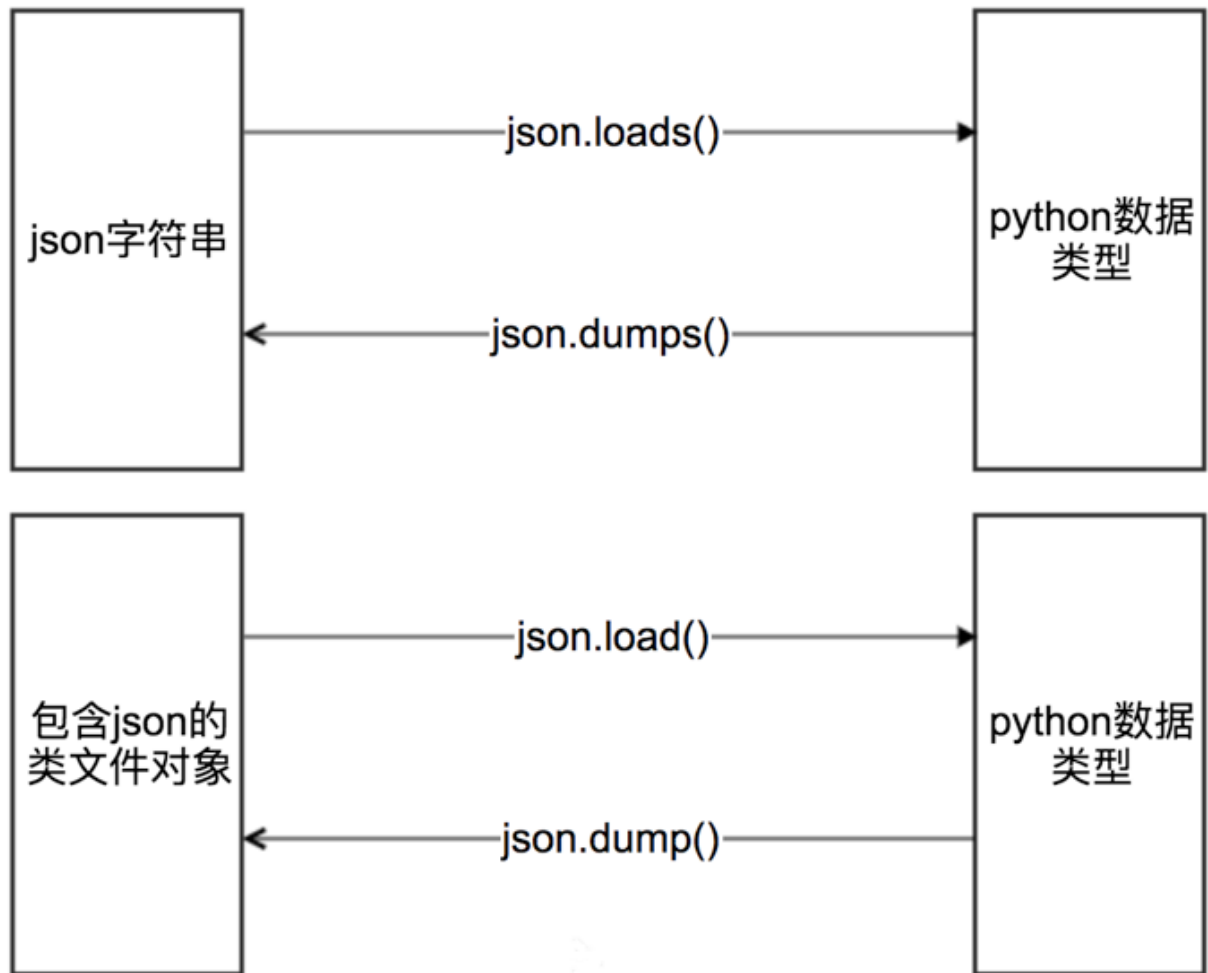
json简单说就是javascript中的对象和数组，所以这两种结构就是对象和数组两种结构，通过这两种结构可以表示各种复杂的结构

对象：对象在js中表示为{ }括起来的内容，数据结构为 { key: value, key: value, ... }的键值对的结构，在面向对象的语言中，key为对象的属性，value为对应的属性值，所以很容易理解，取值方法为 对象.key 获取属性值，这个属性值的类型可以是数字、字符串、数组、对象这几种。

数组：数组在js中是中括号[]括起来的内容，数据结构为 ["Python", "javascript", "C++", ...]，取值方式和所有语言中一样，使用索引获取，字段值的类型可以是 数字、字符串、数组、对象几种。

Python中的json模块

json模块提供了四个功能：dumps、dump、loads、load，用于字符串 和 python数据类型间进行转换。



简记: 只要和str类型有关的就加s, 由str->python对象就是导入, 用loads, 由python对象非->str就是导出, 用dumps

<1> json.loads()

把json格式字符串解码转换成Python对象 从json到python的类型转化对照如下:

```
1 import json
2
3 str_li = "[1, 2, 3, 4]"
4
5 str_dict = '{"name": "JunGe", "city": "Changsha"}'
6
7 json.loads(str_li)
8
9 json.loads(str_dict)
```

<2> json.dumps()

实现python类型转化为json字符串，返回一个str对象 把一个Python对象编码转换成json字符串

```
1 import json
2
3 list_obj = [1, 2, 3, 4]
4
5 dict_obj = {"name": "JunGe", "city": "Changsha"}
6
7 json.dumps(list_obj)
8
9 json.dumps(dict_obj)
```

<3> json.dump()

将Python内置类型序列化为json对象后写入文件

```
1 import json
2
3 list_obj = [{"city": "Changsha"}, {"name": "JunGe"}]
4 json.dump(list_obj, open("list_obj.json", "w"),
5           ensure_ascii=False)
6
7 dict_obj = {"city": "Changsha", "name": "JunGe"}
8 json.dump(dict_obj, open("dict_obj.json", "w"),
9           ensure_ascii=False)
```

<4> json.load()

读取文件中json形式的字符串元素 转化成python类型

```
1 import json
2
3 list_obj = json.load(open("list_obj.json"))
4 print(list_obj)
5
6 dict_obj = json.load(open("dict_obj.json"))
7 print(dict_obj)
```