

xpath解析

XPath是一门在 XML 文档中查找信息的语言. XPath可用来在 XML 文档中对元素和属性进行遍历. 而我们熟知的HTML恰巧属于XML的一个子集. 所以完全可以用xpath去查找html中的内容.

首先, 先了解几个概念.

```
1 <book>
2     <id>1</id>
3     <name>野花遍地香</name>
4     <price>1.23</price>
5     <author>
6         <nick>周大强</nick>
7         <nick>周芷若</nick>
8     </author>
9 </book>
```

在上述html中,

1. book, id, name, price....都被称为节点.
2. Id, name, price, author被称为book的子节点
3. book被称为id, name, price, author的父节点

4. id, name, price,author被称为同胞节点

OK~ 有了这些基础知识后, 我们就可以开始了解xpath的基本语法了
在python中想要使用xpath, 需要安装lxml模块.

```
1 pip install lxml
```

用法:

1. 将要解析的html内容构造出etree对象.
2. 使用etree对象的xpath()方法配合xpath表达式来完成对数据的提取

```
1 from lxml import etree
2
3 html = """
4 <book>
5     <id>1</id>
6     <name>野花遍地香</name>
7     <price>1.23</price>
8     <nick>臭豆腐</nick>
9     <author>
10         <nick id="10086">周大强</nick>
11         <nick id="10010">周芷若</nick>
12         <nick class="joy">周杰伦</nick>
13         <nick class="jolin">蔡依林</nick>
```

```

14         <div>
15             <nick>惹了</nick>
16         </div>
17     </author>
18
19     <partner>
20         <nick id="ppc">胖胖陈</nick>
21         <nick id="ppbc">胖胖不陈</nick>
22     </partner>
23 </book>
24 """
25
26 et = etree.XML(html)
27 # 根据节点进行搜索
28 # result = et.xpath("/book")
29 # result = et.xpath("/book/id") # /在开头表示文档最
    开始，/在中间表示儿子
30 # result = et.xpath("/book//nick") # //表示后代
31 result = et.xpath("/book/*/nick") # *表示通配符
32
33 print(result)
34

```

xpath如何提取属性信息. 我们上一段真实的HTML来给各位讲解一下

准备HTML:

```
1 <!DOCTYPE html>
2 <html lang="en">
3 <head>
4     <meta charset="UTF-8" />
5     <title>Title</title>
6 </head>
7 <body>
8     <ul>
9         <li><a href="http://www.baidu.com">百度
10        </a></li>
11        <li><a href="http://www.google.com">谷
12        歌</a></li>
13        <li><a href="http://www.sogou.com">搜狗
14        </a></li>
15    </ul>
16    <ol>
17        <li><a href="feiji">飞机</a></li>
18        <li><a href="dapao">大炮</a></li>
19        <li><a href="huoche">火车</a></li>
20    </ol>
21    <div class="job">李嘉诚</div>
22    <div class="common">胡辣汤</div>
23 </body>
24 </html>
```

xpath解析

```
1 from lxml import etree
2
3 tree = etree.parse("1.html")
4 result = tree.xpath("/html/body/ul/li/a/@href")
5 print(result)
6
7 result = tree.xpath("/html/body/ul/li")
8 for li in result:
9     print(li.xpath("./a/@href")) # 局部解析
10
11 result = tree.xpath("//div[@class='job']/text()")
12     # [@class='xxx']属性选取 text()获取文本
13 print(result)
```