湯哥聊Python

# 涛哥 优质资料整理

# 1.什么是Beautiful Soup?

Beautiful Soup提供了一种简单而直观的方式来遍历文档的标签树,查找特定的标签或数据,修改文档的结构,以及处理不同的标签和属性。它能够处理不同类型的解析器,包括Python的标准解析器(html.parser)、lxml解析器和html5lib解析器,以适应不同的解析需求。

# 2.Beautiful Soup的用途

- 1. 网页数据抓取: Beautiful Soup用于从网页上抓取数据,包括文本、图片、链接、表格、列表等。
- 2. 信息检索: 它能够帮助搜索和检索特定的标签、文本或属性,以提取感兴趣的信息。
- 3. 数据提取和处理: Beautiful Soup可以解析HTML或XML文档,并从中提取数据,使其易于处理、分析和转换。
- 4. 网页结构分析: 用于分析网页的结构,帮助理解网页的组织和标签之间的层次关系。
- 5. 数据清洗和转换:可以用于清洗和转换从网页抓取的数据,以便进一步的分析和存储。
- 6. 网页爬虫开发:Beautiful Soup是构建网页爬虫的重要工具之一,可以用于创建自动化的数据抓取工具。
- 7. 网站测试和验证:在Web开发中,Beautiful Soup有助于验证和测试网页的结构和内容。

# 3.Beautiful Soup的优点

- 1. 简单而直观: Beautiful Soup提供了简单直观的API,使其易于学习和使用,尤其适合初学者。
- 2. 强大的解析能力:它支持不同类型的解析器,包括Python的标准解析器、lxml和html5lib,因此可以解析各种类型的HTML和XML文档。
- 3. 灵活性: Beautiful Soup允许根据需要轻松地遍历、查找、修改和处理文档的各个部分。
- 4. 容错性:即使在面对损坏或不规范的HTML文档时,Beautiful Soup也能够容忍并尽可能地解析其中的信息。
- 5. Python生态系统的一部分:作为Python库,Beautiful Soup能够与其他Python库和框架良好地集成,如请求库(requests)和数据分析库(Pandas)。

# 二、安装和准备

# 1.安装Beautiful Soup

要安装Beautiful Soup,需要使用Python的包管理工具pip。Beautiful Soup有两个主要的版本:Beautiful Soup 4 (简称BS4) 和 Beautiful Soup 3。通常,建议使用Beautiful Soup 4,因为它更加现代化和功能丰富。

以下是安装Beautiful Soup 4的步骤:

- 1. 打开命令行终端或命令提示符(Windows)。
- 2. 使用以下命令来安装Beautiful Soup 4:
  - 1 pip install beautifulsoup4

这将使用pip从Python包索引中下载Beautiful Soup 4并安装它。如果Python环境有多个版本(例如Python 2和Python 3),请确保使用适当版本的pip。

3. 安装成功后,现在可以在Python中使用Beautiful Soup 4了。

在安装Beautiful Soup后,还需要安装一个HTML解析器,以便Beautiful Soup能够解析HTML文档。 Beautiful Soup支持多个解析器,包括:

- Python的内置解析器: 'html.parser'
- lxml解析器: 'lxml'
- html5lib解析器: 'html5lib'

可以根据自己的需求选择其中一个解析器,但通常来说,lxml是一个性能较好的选择。要安装lxml解析器,可以运行以下命令:

1 pip install lxml

一旦Beautiful Soup和解析器都安装好,就可以在Python中开始使用Beautiful Soup来解析和操作HTML文档了。首先,导入Beautiful Soup库,然后创建Beautiful Soup对象,传入要解析的HTML文档和解析器类型,如下所示:

- 1 from bs4 import BeautifulSoup
  2
- 3 # 创建Beautiful Soup对象
- 4 soup = BeautifulSoup(html, 'lxml')

其中, html 是要解析的HTML文档字符串。

现在,已经成功安装Beautiful Soup并准备好开始使用它来处理HTML文档了。

# 2.导入Beautiful Soup库

要导入Beautiful Soup库,使用以下语句:

```
1 from bs4 import BeautifulSoup
```

这将从Beautiful Soup库中导入Beautiful Soup类,使可以在Python脚本或程序中使用Beautiful Soup来解析和处理HTML或XML文档。

一旦导入了Beautiful Soup库,就可以创建Beautiful Soup对象并开始使用它来解析和操作文档。 以下是一个简单的导入Beautiful Soup库和创建Beautiful Soup对象的示例:

```
1 from bs4 import BeautifulSoup
 2
 3 # 假设html是要解析的HTML文档
4 html = """
 5 <html>
 6 <head>
      <title>Sample HTML Page</title
 8 </head>
9 <body>
       Hello, Beautiful Soup!
10
11 </body>
12 </html>
13 """
14
15 # 创建Beautiful Soup对象
16 soup = BeautifulSoup(html, 'html.parser')
17
    现在可以使用soup来解析和操作HTML文档
18
```

在上面的示例中,首先导入Beautiful Soup库,然后创建了一个Beautiful Soup对象,传入了要解析的 HTML文档字符串和解析器类型(这里使用的是内置解析器'html.parser')。一旦创建了Beautiful Soup对象,就可以使用它来查找、遍历和修改HTML文档中的标签和内容。

这是开始使用Beautiful Soup库的第一步,后续可以根据需要使用它来执行更多复杂的任务,如标签查找、数据提取、数据处理等。

# 3.下载HTML页面

要下载HTML页面,可以使用Python的请求库(例如 requests )来获取页面的内容。

以下是一些示例代码,演示如何使用 requests 库下载HTML页面:

```
1 import requests
3 # 定义要下载的页面的URL
4 url = "https://www.example.com" # 请将此URL替换为要下载的页面URL
6 # 发送HTTP GET请求以获取页面内容
7 response = requests.get(url)
8
9 # 检查请求是否成功
10 if response.status_code == 200:
      # 获取页面内容
11
      html_content = response.text
12
      # 此时,html content包含了所请求页面的HTML内容
13
14 else:
    print("Failed to retrieve the webpage. Status code:", response.status_code)
15
```

在这个示例中,首先导入 requests 库,然后定义了要下载的页面的URL(请替换为感兴趣的页面的URL)。接下来,使用 requests.get(url) 发送HTTP GET请求,以获取页面的内容。如果请求成功(HTTP状态码为200),则 html\_content 变量中将包含页面的HTML内容。

可以将这段代码嵌入到项目中,以便下载特定页面的HTML内容,然后可以将下载的HTML内容传递给 Beautiful Soup对象来进行解析和处理。这是构建网络爬虫或进行网页数据提取的第一步。

# 三、Beautiful Soup基础

# 1.创建Beautiful Soup对象

创建Beautiful Soup对象,需要导入Beautiful Soup库并使用其构造函数来解析HTML或XML文档。 以下是创建Beautiful Soup对象的一般步骤:

1. 导入Beautiful Soup库:

```
1 from bs4 import BeautifulSoup
```

2. 定义要解析的HTML或XML文档。这可以是从Web下载的页面内容,也可以是本地文件中的文档内容。

3. 使用Beautiful Soup构造函数创建Beautiful Soup对象,同时指定解析器的类型。通常,可以选择内置的'html.parser'解析器或其他解析器(如'lxml'或'html5lib')。这里以内置解析器为例:

现在,soup 是一个Beautiful Soup对象,可以用于遍历、查找和操作HTML文档中的标签和内容。如果要解析一个网站上的实际页面,首先需要使用 requests 库或其他HTTP库下载页面内容,然后将该内容传递给Beautiful Soup构造函数,

#### 如下所示:

```
1 import requests
2 from bs4 import BeautifulSoup
3
4 # 定义要下载的页面的URL
5 url = "https://www.example.com" # 请将此URL替换为要下载的页面URL
  # 发送HTTP GET请求以获取页面内容
   response = requests.get(url)
8
9
     检查请求是否成功
10
  if response.status_code == 200:
       # 获取页面内容
12
13
      html_content = response.text
14
      # 创建Beautiful Soup对象
15
16
      soup = BeautifulSoup(html_content, 'html.parser')
17
      # 现在可以使用soup来解析和操作HTML文档
18
19 else:
      print("Failed to retrieve the webpage. Status code:", response.status_code)
20
```

这样,就创建了一个Beautiful Soup对象,可以使用它来进行HTML文档的解析和操作。从此步骤开始,可以使用Beautiful Soup提供的方法来查找标签、提取数据、遍历文档,以及执行其他操作。

# 2.Beautiful Soup对象的基本属性

Beautiful Soup对象具有一些基本属性,这些属性允许访问文档的不同部分和信息。

以下是一些Beautiful Soup对象的基本属性:

1. **.title**: 获取HTML文档中的标题标签。

```
1 title_tag = soup.title
2 print(title_tag)
```

2. head: 获取HTML文档中的头部标签。

```
1 head_tag = soup.head
2 print(head_tag)
```

3. **.body**: 获取HTML文档中的主体标签。

```
1 body_tag = soup.body
2 print(body_tag)
```

4. **prettify()**:格式化整个文档的内容,使其更易于阅读。这个方法通常用于调试和输出HTML内容。

```
1 pretty_html = soup.prettify()
2 print(pretty_html)
```

这些属性允许访问Beautiful Soup对象表示的HTML或XML文档的不同部分。可以根据需要使用这些属性来访问文档的标题、头部、主体或对整个文档进行格式化,以便更清晰地查看文档的结构和内容。

以下是一个示例,演示如何使用Beautiful Soup对象的基本属性:

```
1 from bs4 import BeautifulSoup
 2
 3 html = """
 4 <html>
 5 <head>
 6
      <title>Sample HTML Page</title>
 7 </head>
 8 <body>
       Hello, Beautiful Soup!
10 </body>
11 </html>
12 """
13
14 soup = BeautifulSoup(html, 'html.parser')
15
16 # 获取标题标签
17 title_tag = soup.title
18 print("Title Tag:", title_tag)
19
20 # 获取头部标签
21 head_tag = soup.head
22 print("Head Tag:", head_tag)
23
24 # 获取主体标签
25 body_tag = soup.body
26 print("Body Tag:", body_tag)
27
28 # 格式化整个文档
29 pretty_html = soup.prettify()
30 print("Pretty HTML:")
31 print(pretty_html)
```

上述代码演示了如何使用Beautiful Soup对象的基本属性来访问HTML文档的不同部分,并以格式化的方式输出整个文档。

# 3.Beautiful Soup对象的基本方法

Beautiful Soup对象提供了许多基本方法,用于解析、查找、遍历和操作HTML或XML文档中的标签和内容。

以下是一些Beautiful Soup对象的基本方法:

1. .find(name, attrs, recursive, string, \*\*kwargs): 查找并返回第一个匹配条件的标签。

```
1 tag = soup.find('a', class_='link')
```

2. .find\_all(name, attrs, recursive, string, limit, \*\*kwargs): 查找并返回 所有匹配条件的标签,以列表形式返回。 1 tags = soup.find\_all('a', class\_='link') 3. .select(css\_selector):使用CSS选择器语法查找标签。 1 tags = soup.select('div > p') 4. .get\_text(): 获取标签内所有文本内容,包括嵌套标签。 1 text = tag.get\_text() 5. .get(name, default): 获取标签的属性值。 1 href = tag.get('href', 'No link') 6. .has\_attr(name): 检查标签是否包含特定属性。 1 has\_href = tag.has\_attr('href') 7. .name: 获取标签的名称。 1 tag\_name = tag.name 8. .string: 获取标签内的文本内容,不包括嵌套标签。 1 tag\_text = tag.string

9. .find\_parent(name, attrs, recursive, string, \*\*kwargs): 查找父标签。

```
1 parent_tag = tag.find_parent('div', class_='container')
```

10. find\_next\_sibling(name, attrs, recursive, string, \*\*kwargs): 查找下一个兄弟标签。

```
1 next_sibling = tag.find_next_sibling('p')
```

11. .find\_all\_next(name, attrs, recursive, string, \*\*kwargs): 查找之后的所有匹配标签。

```
1 next_tags = tag.find_all_next('a', class_='link')
```

12. .find\_previous\_sibling(name, attrs, recursive, string, \*\*kwargs): 查找上一个兄弟标签。

```
1 prev_sibling = tag.find_previous_sibling('p')
```

13. .find\_all\_previous(name, attrs, recursive, string, \*\*kwargs): 查找之前的所有匹配标签。

```
1 prev_tags = tag.find_all_previous('a', class_='link')
```

这些方法允许在Beautiful Soup对象上执行各种操作,例如查找特定标签、提取数据、遍历文档、检查属性和属性值,以及执行其他操作。可以根据需要使用这些方法来满足数据提取和处理需求。

# 四、解析HTML页面

# 1.解析HTML页面的不同方式

Beautiful Soup可以使用不同类型的解析器来解析HTML页面。解析器是Beautiful Soup用来解析HTML文档的引擎,不同的解析器有不同的优点和适用场景。

以下是解析HTML页面的不同方式:

### 1. 使用Python的标准解析器('html.parser'):

这是Beautiful Soup的内置解析器,通常是最快的解析器,而且在大多数情况下都能很好地工作。它是Beautiful Soup的默认解析器。

```
1 soup = BeautifulSoup(html, 'html.parser')
```

### 2. 使用lxml解析器:

lxml是一个第三方解析器,它具有出色的性能和容错性,可以处理大型、复杂的HTML文档。

```
1 soup = BeautifulSoup(html, 'lxml')
```

### 3. 使用html5lib解析器:

html5lib是一个纯Python实现的解析器,它能够处理不规范的HTML文档,并具有很好的容错性。它通常更慢一些,但在处理一些特殊情况下非常有用,例如处理不规范的HTML。

```
1 soup = BeautifulSoup(html, 'html5lib')
```

选择解析器的方式取决于需求和HTML文档的特性。通常情况下,'html.parser'是一个不错的默认选择,因为它性能良好,并在大多数情况下都能正常工作。如果需要更好的性能或需要处理复杂的文档,可以考虑使用lxml解析器。而html5lib解析器则适用于处理不规范的文档。

要使用不同的解析器,只需将解析器的名称作为第二个参数传递给Beautiful Soup构造函数即可。无论使用哪个解析器,Beautiful Soup的API都是一样的,因此可以在不同的解析器之间无缝切换。

# 2.标签的基本结构

HTML标签是构成HTML文档结构的基本元素,具有一定的结构。HTML标签通常由开始标签和结束标签组成,有些标签也可以是自闭合的。

以下是HTML标签的基本结构:

1. **开始标签**: 开始标签由尖括号 < 和标签名称组成,标签名称通常是HTML元素的名称。开始标签用于定义标签的起始点。

例如, 是一个 标签的开始标签。

2. **结束标签**:结束标签也由尖括号 < 、斜杠 / 和标签名称组成,标签名称与开始标签中的名称相同。 结束标签用于定义标签的结束点。

例如, 是一个 标签的结束标签。

3. **自闭合标签**:有一些HTML标签是自闭合的,它们不需要结束标签。自闭合标签通常以 < tag /> 的形式表示,其中 tag 是标签名称。

例如, <img src="image.jpg" /> 是一个自闭合的 <img> 标签,用于显示图像。

4. **属性**:标签可以包含属性,属性位于开始标签中,通常以键值对的形式表示。属性用于提供有关标签的额外信息。

例如, <a href="https://www.example.com"> 中的 href 是 <a> 标签的属性,它指定链接的URL。

5. **标签内容**:在开始标签和结束标签之间的部分称为标签的内容。内容可以是文本、其他嵌套标签或 其他元素。

例如, Hello, World! 中的 Hello, World! 是 标签的内容。

总结一下, HTML标签的基本结构如下:

• 开始标签: <tag>

• 结束标签: </tag>

• 自闭合标签: <tag />

• 属性: attribute="value"

• 标签内容: Content

HTML标签是HTML文档的构建块,通过组合不同的标签和属性,可以创建丰富的Web页面。在使用 Beautiful Soup或进行HTML文档操作时,了解标签的基本结构非常重要。

### 3.寻找标签

Beautiful Soup提供了多种方法来查找和定位HTML文档中的标签。可以按照标签名、类名、ID、属性、文本内容以及嵌套标签等不同的方式来查找标签。

以下是示例:

### 1. 按标签名查找标签:

使用 find 方法或 find\_all 方法按标签名查找标签。 find 会返回第一个匹配的标签,而 find\_all 会返回所有匹配的标签。

```
2 link = soup.find('a')
3
4 # 查找所有标签
5 paragraphs = soup.find_all('p')
```

### 2. 按类名查找标签:

使用 find 方法或 find\_all 方法按类名查找标签。将类名作为参数传递给 class\_ 参数。

```
1 # 查找第一个具有类名为 "example" 的<div>标签
2 div = soup.find('div', class_='example')
3
4 # 查找所有具有类名为 "highlight" 的<span>标签
5 spans = soup.find_all('span', class_='highlight')
```

### 3. 按ID查找标签:

使用 find 方法按ID查找标签,将ID名作为参数传递给 id 参数。

```
1 # 查找具有ID为 "header" 的<div>标签
2 header = soup.find('div', id='header')
```

#### 4. 按属性查找标签:

使用 find 方法或 find\_all 方法按属性查找标签。将属性名和属性值作为参数传递给字典。

```
1 # 查找所有<a>标签,其中href属性包含 "example.com"
2 links = soup.find_all('a', {'href': 'example.com'})
3
4 # 查找所有<img>标签,其中alt属性包含 "logo"
5 logos = soup.find_all('img', {'alt': 'logo'})
```

### 5. 按文本内容查找标签:

使用 find 方法或 find\_all 方法按文本内容查找标签。将要匹配的文本作为参数传递。

```
1 # 查找第一个包含文本 "Hello, World!" 的标签
2 paragraph = soup.find('p', text='Hello, World!')
3
4 # 查找所有包含文本 "Example" 的<span>标签
5 examples = soup.find_all('span', text='Example')
```

### 6. 查找嵌套标签:

使用标签对象的属性来查找嵌套标签。例如,如果有一个包含父标签的变量,可以使用点号来查找子标签。

- 1 # 假设div tag是一个包含<div>标签的Beautiful Soup对象
- 2 # 查找在<div>标签内的标签
- 3 inner\_paragraph = div\_tag.find('p')

这些方法允许根据不同的标签特征来查找和定位HTML文档中的标签。可以根据具体的需求选择合适的 查找方式。

# 五、遍历HTML文档

# 1.遍历父标签和子标签

Beautiful Soup允许轻松地遍历父标签和子标签,以访问HTML文档中的不同层次的标签。可以使用 Beautiful Soup提供的属性和方法来实现这些遍历操作。

以下是如何遍历父标签和子标签的示例:

### 遍历子标签:

可以使用Beautiful Soup的 .contents 属性或 .children 属性来遍历子标签。这允许访问父标签内的直接子标签。

```
1 from bs4 import BeautifulSoup
2
3 html = """
   <html>
5
       <body>
        <div>
               Paragraph 1
               Paragraph 2
           </div>
       </body>
11 </html>
   11 11 11
12
13
14 soup = BeautifulSoup(html, 'html.parser')
15
```

```
16 # 遍历<body>标签内的子标签
17 body_tag = soup.body
18 for child in body_tag.contents:
19 print(child)
```

上述代码遍历了 <body> 标签内的子标签,并打印了每个子标签。

### 遍历父标签:

可以使用 .parent 属性来遍历父标签。这允许访问当前标签的直接父标签。

```
1 # 假设p_tag是一个标签的Beautiful Soup对象
2 # 遍历标签的父标签
3 parent_tag = p_tag.parent
4 print(parent_tag)
```

上述代码演示了如何遍历 标签的父标签。

### 遍历所有子标签:

如果想遍历父标签下的所有子标签,包括子标签的子标签,可以使用 .descendants 属性。

```
1 # 遍历<body>标签内的所有子标签,包括嵌套的标签
2 for tag in body_tag.descendants:
3 print(tag)
```

上述代码会遍历 <body> 标签内的所有子标签,包括嵌套的标签。

这些方法使能够灵活地遍历HTML文档的标签层次结构,以便访问和操作其中的内容。可以根据需要选择适当的遍历方法来定位和处理特定的标签。

# 2.使用循环遍历标签

可以使用循环来遍历HTML文档中的标签,以查找、处理或操作多个标签。

以下是如何使用循环遍历标签的示例:

```
8
               Paragraph 2
           </div>
 9
           <div>
10
               <a href="https://www.example.com">Link 1</a>
11
               <a href="https://www.example.org">Link 2</a>
12
13
           </div>
       </body>
14
15 </html>
  111111
16
17
18 soup = BeautifulSoup(html, 'html.parser')
19
20 # 遍历所有 <a>标签
21 for a_tag in soup.find_all('a'):
       print("Link Text:", a_tag.text)
22
23
       print("Link URL:", a_tag.get('href'))
       print()
24
25
26 # 遍历所有 < div > 标签
27 for div_tag in soup.find_all('div'):
       print("Contents of <div>:", div_tag.text)
28
       print()
29
```

在上述示例中,使用 find\_all 方法和循环分别遍历了所有 <a> 标签和所有 <div> 标签。在每次 迭代中,可以访问标签的属性、文本内容以及执行其他操作。

使用循环遍历标签时,可以根据标签的名称、类名、ID、属性等条件来筛选和处理特定的标签。这允许灵活地处理HTML文档中的多个标签,并提取需要的信息。

# 3.获取标签的属性和文本内容

要获取标签的属性和文本内容,可以使用Beautiful Soup提供的属性和方法。

以下是如何获取标签的属性和文本内容的示例:

### 获取标签的属性:

可以使用 .get(attribute\_name) 方法来获取标签的属性值。在括号中传递属性的名称,方法会返回属性的值。

```
1 from bs4 import BeautifulSoup
2
3 html = """
4 <a href="https://www.example.com" class="link">Example Link</a>
5 """
6
7 soup = BeautifulSoup(html, 'html.parser')
```

```
8
9 # 获取<a>标签的href属性
10 href = soup.a.get('href')
11 print("href:", href)
12
13 # 获取<a>标签的class属性
14 class_attr = soup.a.get('class')
15 print("class:", class_attr)
```

上述代码演示了如何获取 <a> 标签的 href 属性和 class 属性的值。

### 获取标签的文本内容:

可以使用 text 属性来获取标签内的文本内容。这将返回标签内所有文本的合并字符串,包括嵌套标签的文本。

```
1 # 获取<a>标签内的文本内容
2 text_content = soup.a.text
3 print("Text Content:", text_content)
```

上述代码演示了如何获取 <a> 标签内的文本内容,它会返回标签内的文本 "Example Link"。

如果只希望获取标签内的直接文本内容(不包括嵌套标签的文本),可以使用 .string 属性:

```
1 # 获取<a>标签内的直接文本内容
2 direct_text_content = soup.a.string
3 print("Direct Text Content:", direct_text_content)
```

.string 属性只返回标签内的直接文本内容,不包括嵌套标签的文本。

使用这些方法,可以轻松地获取标签的属性和文本内容,以进一步处理和操作HTML文档中的数据。

# 六、修改HTML文档

# 1.修改标签的属性

要修改标签的属性,可以使用Beautiful Soup来选择目标标签,然后使用标签对象的属性来更改属性值。

以下是如何修改标签的属性的示例:

```
1 from bs4 import BeautifulSoup
```

```
3 html = """
4 <a href="https://www.example.com" class="link">Example Link</a>
6
7 soup = BeautifulSoup(html, 'html.parser')
8
9 # 选择目标<a>标签
10 a_tag = soup.a
11
12 # 修改href属性的值
13 a_tag['href'] = 'https://www.newlink.com'
14
15 # 修改class属性的值,可以使用空格分隔多个类名
16 a_tag['class'] = 'new-class1 new-class2'
17
18 # 打印修改后的标签
19 print(a_tag)
```

上述代码演示了如何选择 <a> 标签,然后分别修改其 href 和 class 属性的值。可以通过简单地将新的值分配给标签对象的属性来修改属性值。

#### 注意:

- 1. 在修改属性时,请确保属性名称(如 'href' 和 'class' )是正确的。
- 2. 如果目标标签中不存在要修改的属性,将添加新属性。
- 3. 如果要将属性值移除,可以使用 del 关键字。例如, del a\_tag['class'] 将删除 <a> 标 签的 class 属性。

通过这种方式,可以轻松地更改标签的属性,以满足特定需求。

# 2.修改标签的文本内容

要修改标签的文本内容,可以使用Beautiful Soup来选择目标标签,然后使用标签对象的 .string 属性或 .replace\_with() 方法来更改文本内容。

以下是如何修改标签的文本内容的示例:

### 使用 .string 属性修改文本内容:

```
1 from bs4 import BeautifulSoup
2
3 html = """
```

```
4 This is the original text.
5 """
6
7 soup = BeautifulSoup(html, 'html.parser')
8
9 # 选择目标标签
10 p_tag = soup.p
11
12 # 修改文本内容
13 p_tag.string = "This is the new text."
14
15 # 打印修改后的标签
16 print(p_tag)
```

上述代码演示了如何选择 标签,然后使用 .string 属性将文本内容修改为新的值。

### 使用 .replace\_with() 方法修改文本内容:

```
1 from bs4 import BeautifulSoup
3 html = """
4 This is the original text.
7 soup = BeautifulSoup(html, 'html.parser')
8
9 # 选择目标标签
10 p_tag = soup.p
11
12 # 创建新文本
13 new_text = "This is the new text."
14
15 # 使用replace_with()方法替换文本内容
16 p_tag.string.replace_with(new_text)
17
18 # 打印修改后的标签
19 print(p_tag)
```

在上述示例中,首先选择了 标签,然后使用 .replace\_with() 方法将文本内容替换为新的文本。

# 3.插入新标签

要在HTML文档中插入新标签,可以使用Beautiful Soup提供的方法来创建新标签,然后将新标签插入 到特定的位置。

以下是如何插入新标签的示例:

### 插入新标签到指定位置:

```
1 from bs4 import BeautifulSoup
2
3 html = """
4 <div>
     Existing paragraph
6 </div>
7 """
8
9 soup = BeautifulSoup(html, 'html.parser')
10
11 # 创建一个新的<a>标签
12 new_a_tag = soup.new_tag('a')
13
14 # 设置新标签的属性
15 new_a_tag['href'] = 'https://www.example.com
16 new_a_tag.string = 'New Link'
17
18 # 找到要插入的位置,例如在<di→标签内插入
19 div_tag = soup.div
20
21 # 使用.append()方法将新标签插入到指定位置
22 div_tag.append(new_a_tag)
23
24 # 打印修改后的HTM
25 print(soup)
```

在上述示例中,首先创建了一个新的〈a〉标签,设置了其属性和文本内容。然后,选择要插入新标签的位置,这里选择的是〈div〉标签,最后使用 .append() 方法将新标签插入到指定位置。

### 插入新标签后的HTML输出:

```
1 <div>
2 Existing paragraph
3 <a href="https://www.example.com">New Link</a>
4 </div>
```

还可以使用 .insert() 方法、 .insert\_before() 方法、 .insert\_after() 方法等,根据需要选择合适的插入方式。这些方法允许在HTML文档中插入新标签,并以不同的方式控制插入位置。

### 4.删除标签

要删除标签,可以使用Beautiful Soup提供的方法,如 .decompose() 或 .extract() 来删除目标标签。

以下是如何删除标签的示例:

# 使用 .decompose() 方法删除标签:

```
1 from bs4 import BeautifulSoup
3 html = """
4 <div>
     Paragraph 1
      Paragraph 2
7 </div>
  0.00
8
9
10 soup = BeautifulSoup(html, 'html.parser')
11
12 # 选择要删除的标签
13 p_tag = soup.find('p')
14
15 # 使用 .decompose() 方法删除标签
16 p_tag.decompose()
17
18 # 打印修改居的HTML
19 print(soup)
```

在上述示例中,首先选择了要删除的 标签,然后使用 .decompose() 方法将其从文档中删除。

### 使用 .extract() 方法删除标签:

在上述示例中,同样选择了要删除的 标签,然后使用 .extract() 方法将其从文档中删除。

# 七、实际应用

## 1.爬取网页数据

要爬取网页数据,可以使用Python的Beautiful Soup库与HTTP请求库(如requests)一起工作。

以下是爬取网页数据的一般步骤:

1. 导入必要的库:

```
1 from bs4 import BeautifulSoup
2 import requests
```

2. 发送HTTP请求并获取网页内容:

使用HTTP请求库发送GET请求,获取网页的HTML内容。

```
1 url = 'https://www.example.com' # 替换为要爬取的网页地址
2 response = requests.get(url)
3 html = response.text
```

3. 创建Beautiful Soup对象:

使用Beautiful Soup库解析HTML内容,创建Beautiful Soup对象。

```
1 soup = BeautifulSoup(html, 'html.parser')
```

#### 4. 查找和提取数据:

使用Beautiful Soup的方法来查找和提取需要的数据。可以使用标签名称、类名、属性、文本内容等条件来定位数据。

```
1 # 示例: 查找所有<a>标签的链接
2 links = soup.find_all('a')
3 for link in links:
4 print(link['href'])
```

### 5. 数据处理和存储:

对提取的数据进行处理或存储,例如保存到文件、数据库或进行进一步分析。

下面是一个完整的示例,演示了如何爬取网页数据:

```
1 from bs4 import BeautifulSoup
2 import requests
3
4 # 发送HTTP请求并获取网页内容
5 url = 'https://www.example.com'
6 response = requests.get(url)
7 html = response.text
8
9 # 创建Beautiful Soup对象
10 soup = BeautifulSoup(html, 'html.parser')
11
12 # 查找和提取数据 这里查找所有<a>标签的链接
13 links = soup.find_all('a')
14 for link in links:
15 print(link['href'])
```

上述示例爬取了指定网页的所有链接,并打印了它们的URL。可以根据需要修改代码以满足网页数据爬取需求。请注意,合法和道德的爬取是非常重要的,确保有权爬取和使用网页上的数据,并遵守网站的使用政策和法律法规。

# 2.数据提取与处理

在网页爬虫中,数据提取和处理是非常重要的一步。Beautiful Soup提供了多种方法来查找和提取数据,而Python的数据处理库(如Pandas)可以用来处理和分析数据。

#### 以下是数据提取和处理的一般步骤:

### 1. 查找和提取数据:

使用Beautiful Soup的查找方法(如 .find() 、 .find\_all() 、 .select() 等)来定位网页中的数据。可以根据标签名称、类名、属性、文本内容等条件来定位数据。

```
1 # 示例: 提取网页中的新闻标题和链接
2 news_titles = soup.select('.news-title')
3 news_links = soup.select('.news-link')
```

### 2. 数据结构化:

将提取的数据结构化为适合进一步处理的数据结构,如列表、字典或Pandas的DataFrame。根据需求,可以使用循环或列表推导来创建数据结构。

```
1 # 示例: 创建包含新闻标题和链接的字典列表
2 news_data = []
3 for title, link in zip(news_titles, news_links):
4     news_data.append({'title': title.text, 'link': link['href']})
```

### 3. 数据处理:

对提取的数据进行处理,可以进行数据清洗、格式转换、计算等操作。Python的数据处理库(如 Pandas)可以在这一步中派上用场。

```
1 import pandas as pd
2
3 # 将数据存储为Pandas DataFrame
4 df = pd.DataFrame(news_data)
```

### 4. 数据分析和可视化:

使用数据处理库来进行数据分析,执行统计分析、可视化等操作。例如,可以使用Matplotlib或 Seaborn来创建图表。

```
1 import matplotlib.pyplot as plt
2
3 # 示例: 绘制新闻标题长度的直方图
4 df['title_length'] = df['title'].apply(len)
```

```
5 plt.hist(df['title_length'], bins=20)
6 plt.xlabel('Title Length')
7 plt.ylabel('Frequency')
8 plt.show()
```

### 5. 数据存储:

将处理后的数据存储到文件、数据库或其他数据存储系统中,以备后续使用。

```
1 # 示例: 将数据保存为CSV文件
2 df.to_csv('news_data.csv', index=False)
```

数据提取和处理是网页爬虫的关键部分,它们决定了能够从网页中获取的信息以及如何使用它们。根据具体的爬取任务,可以执行各种数据提取和处理操作,以满足需求。

### 3.数据存储

将爬取到的数据存储到文件、数据库或其他数据存储系统是网页爬虫的重要一步。具体的存储方式取决于需求和使用场景。

以下是常见的数据存储方式:

### 1. 存储为文本文件(例如CSV):

将数据存储为文本文件是一种常见且简单的方式,适用于小型数据集。可以使用Python的内置库(如 open()),或第三方库(如Pandas)来将数据写入CSV文件。

```
1 import pandas as pd
2
3 # 将数据存储为CSV文件
4 df.to_csv('data.csv', index=False)
```

### 2. 存储为JSON文件:

JSON是一种常用的数据交换格式,适合保存结构化数据。可以使用Python的内置库(如 j son )来将数据保存为JSON文件。

```
1 import json
2
3 # 将数据存储为JSON文件
```

```
4 with open('data.json', 'w') as json_file:
5    json.dump(data, json_file)
```

### 3. 存储到数据库:

如果需要长期保存数据,将数据存储到数据库是一个不错的选择。可以使用关系型数据库(如 SQLite、MySQL、PostgreSQL)或NoSQL数据库(如MongoDB)来存储数据。

```
1 import sqlite3
2
3 # 连接到SQLite数据库
4 conn = sqlite3.connect('mydata.db')
5
6 # 使用Pandas将数据存储到SQLite数据库
7 df.to_sql('mytable', conn, if_exists='replace', index=False)
8
9 # 关闭数据库连接
10 conn.close()
```

### 4. 使用对象存储服务:

对于大型数据集或需要云存储的情况,可以考虑使用云服务提供商(如Amazon S3、Google Cloud Storage)提供的对象存储服务。

### 5. 存储到内存数据库:

如果只需要临时存储数据,可以考虑使用内存数据库(如SQLite的内存模式)。

```
1 import sqlite3
2
3 # 连接到SQLite内存数据库
4 conn = sqlite3.connect(':memory:')
5
6 # 使用Pandas将数据存储到内存数据库
7 df.to_sql('mytable', conn, if_exists='replace', index=False)
```

选择适当的数据存储方式取决于需求,包括数据的大小、持久性、访问需求以及数据后续的处理和分析。在选择存储方式时,确保遵循数据隐私和法规,以及服务提供商的使用政策。

# 八、示例项目

### 1.创建一个简单的网页爬虫

创建一个简单的网页爬虫需要使用Python的Beautiful Soup库和HTTP请求库(如requests)。 以下是一个简单的示例,演示如何爬取一个网页上的链接:

```
1 from bs4 import BeautifulSoup
2 import requests
4 # 发送HTTP请求并获取网页内容
5 url = 'https://www.example.com' # 替换为要爬取的网页地址
6 response = requests.get(url)
7 html = response.text
8
9 # 创建Beautiful Soup对象
10 soup = BeautifulSoup(html, 'html.parser')
11
12 # 查找所有<a>标签的链接
13 links = soup.find_all('a')
14
15 # 打印链接
16 for link in links:
     print(link['href'])
17
```

在上述示例中,首先发送HTTP GET请求获取网页内容,然后使用Beautiful Soup解析HTML。随后,使用 .find\_all() 方法查找所有 <a> 标签,然后打印它们的链接。

# 2.提取特定网页的数据

要提取特定网页的数据,需要首先分析目标网页的结构和内容,然后使用Beautiful Soup来定位和提取需要的数据。以下是一个通用的示例,展示如何提取特定网页的数据:

假设要从一个示例网页(https://www.example.com)中提取标题和链接的数据:

```
1 from bs4 import BeautifulSoup
2 import requests
3
4 # 发送HTTP请求并获取网页内容
5 url = 'https://www.example.com' # 替换为要爬取的网页地址
6 response = requests.get(url)
7 html = response.text
8
9 # 创建Beautiful Soup对象
10 soup = BeautifulSoup(html, 'html.parser')
11
```

```
12 # 查找并提取标题和链接
13 titles = []
14 links = []
15
16 # 示例: 假设标题是<h2>标签, 链接在<a>标签中
17 title_tags = soup.find_all('h2')
18 link_tags = soup.find_all('a')
19
20 # 提取标题
21 for tag in title_tags:
      titles.append(tag.text)
22
23
24 # 提取链接
25 for tag in link_tags:
      links.append(tag['href'])
26
27
28 # 打印提取的数据
29 for title, link in zip(titles, links):
      print(f"Title: {title}")
30
      print(f"Link: {link}")
31
32
      print()
```

上述示例中,首先获取网页内容,然后使用Beautiful Soup解析HTML。接着,查找和提取标题(假设它们在〈h2〉标签中)和链接(假设它们在〈a〉标签中)。最后,打印提取的数据。

# 3.将数据保存到文件

将爬取到的数据保存到文件可以使用Python内置的文件操作方法或使用数据处理库,具体方式取决于数据的格式和需求。

以下是示例代码,展示如何将数据保存到文本文件(CSV):

```
1 import csv
2
3 # 假设已经有一个包含数据的列表,例如:
4 data = [
5 {'Title': 'Title 1', 'Link': 'https://www.example.com/1'},
6 {'Title': 'Title 2', 'Link': 'https://www.example.com/2'},
7 {'Title': 'Title 3', 'Link': 'https://www.example.com/3'}
8 ]
9
10 # 指定要保存的文件名
11 filename = 'data.csv'
12
13 # 使用CSV库将数据保存为CSV文件
```

```
14 with open(filename, 'w', newline='') as csv_file:
       fieldnames = ['Title', 'Link']
15
       writer = csv.DictWriter(csv_file, fieldnames=fieldnames)
16
17
      # 写入CSV文件的标题行
18
       writer.writeheader()
19
20
       # 写入数据
21
22
       for item in data:
           writer.writerow(item)
23
24
25 print(f'Data saved to {filename}')
```

在这个示例中,使用Python内置的 csv 库创建了一个CSV文件,然后使用 csv.DictWriter 来写入数据。确保数据格式与示例相匹配,并替换示例数据和字段名。

另外,也可以使用其他数据格式,如JSON、Excel等,根据需要调整保存数据的方式。例如,使用 Pandas库可以轻松地将数据保存为各种格式,如CSV、Excel、JSON等。

```
1 import pandas as pd

2

3 # 将数据存储为CSV文件

4 df = pd.DataFrame(data)

5 df.to_csv('data.csv', index=False)

6

7 # 或将数据存储为Excel文件

8 df.to_excel('data.xlsx', index=False)

9

10 # 或将数据存储为JSON文件

11 df.to_json('data.json', orient='records')
```

根据需求和数据格式,选择适当的保存方式。

# 九、常见问题和注意事项

# 1.处理编码问题

在网页爬虫中,可能会遇到编码问题,尤其是当爬取网页上的文本数据时。处理编码问题是非常重要的,以确保正确地解析和处理网页内容。

以下是一些处理编码问题的建议:

#### 1. 指定网页的编码格式:

在发起HTTP请求时,尽量明确指定网页的编码格式。大多数网页会在响应头部(Response Headers)中提供 Content-Type 字段,其中包含字符编码信息。可以使用 response encoding 属性来指定编码格式。

```
1 response = requests.get(url)
2 response.encoding = 'utf-8' # 指定编码格式为UTF-8
```

### 2. 处理字符编码异常:

使用异常处理来处理字符编码异常。有些网页可能包含无效的字符或编码错误,这可能会导致 UnicodeDecodeError 。在这种情况下,可以使用 try 和 except 块来捕获异常并处理。

```
1 try:
2 html = response.text
3 except UnicodeDecodeError:
4 # 处理编码异常的代码
```

### 3. 使用Beautiful Soup的 . content 属性:

如果仍然遇到编码问题,可以使用Beautiful Soup的 .content 属性来获取原始的字节数据,而不是解码后的文本数据。然后,可以手动指定编码并解码数据。

```
1 content = response.content
2 html = content.decode('utf-8') # 指定编码并解码数据
```

### 4. 使用自动编码检测库:

还可以使用第三方库,如 chardet ,来自动检测网页的编码。这可以帮助确定正确的编码格式。

```
1 import chardet
2
3 charset = chardet.detect(response.content)['encoding']
4 response.encoding = charset
5 html = response.text
```

以上建议中的某一方法通常能够解决大多数编码问题。根据具体情况,可以选择适合方式来处理编码问题,以确保可以正确地解析和处理网页内容。

### 2.处理网站结构变化

当网站的结构发生变化时,需要调整网页爬虫来适应这些变化。网站结构变化可能包括HTML标签的更改、CSS类名的修改、数据位置的移动等。

以下是处理网站结构变化的一些建议:

### 1. 重新分析网页结构:

首先,重新审查网页的结构。使用浏览器的开发者工具(例如Chrome开发者工具)来查看网页的HTML结构、标签和CSS类名。确定哪些部分发生了变化。

### 2. 更新选择器和条件:

根据网站结构的变化,更新Beautiful Soup选择器和条件。可能需要使用不同的标签名称、类名、属性等来查找和提取数据。可以使用 .select() 方法或 .find() 方法来根据新的条件选择元素。

### 3. 处理数据位置的变化:

如果数据的位置发生了变化,需要调整代码以正确找到数据。这可能涉及到深度查找或使用嵌套选择器来定位数据。

### 4. 使用异常处理:

考虑使用异常处理来处理结构变化引起的问题。如果无法找到所需的元素,可以捕获异常并采取相应的措施,如回退到备用选择器或跳过某些数据。

```
1 try:
2 # 尝试查找元素
3 data = soup.find('div', {'class': 'new-class'})
4 except AttributeError:
5 # 处理未找到元素的情况
6 data = None
```

### 5. 定期检查和测试:

定期检查和测试爬虫,以确保它适应网站结构的任何变化。监控网站的变化并及时更新爬虫代码。

#### 6. 遵守网站的使用政策:

← 在更新爬虫时,始终确保遵守网站的使用政策和使用协议。一些网站可能禁止爬虫访问,因此请确保有权爬取和使用网站上的数据。

处理网站结构变化需要不断的观察和调整。根据网站的变化情况,灵活地更新爬虫代码以确保其正常 运行。

# 3.爬取道德和法律考虑事项

在网页爬虫的过程中,需要严格遵守道德和法律规定,以确保爬虫活动是合法和道德的

### 法律考虑事项:

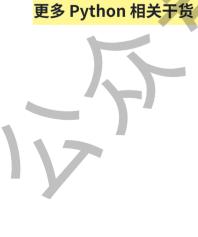
- 1.合法性:确保有权访问和爬取特定网站的内容。某些网站可能有明确的使用政策,可能禁止或限制爬 虫访问。
- **2. 著作权**: 尊重著作权法。只爬取和使用允许爬虫访问和使用的内容,避免未经许可爬取和分发受版 权保护的材料。
- 3. 隐私法:避免爬取和存储用户个人数据,以免侵犯隐私法规。对于敏感信息,必须尊重用户的隐私 权。
- **4. 反垃圾邮件法**:如果计划使用爬取到的数据来发送电子邮件或消息,遵守反垃圾邮件法规、 供取消订阅选项和不发送未经请求的消息。
- **5. 合同和使用政策**: 尊重网站的使用政策和服务条款。某些网站可能要求遵守特定的使用政策,否则 可能会采取法律措施。

### 道德考虑事项:

- 1. 尊重网站的带宽:避免过度频繁的请求,以减轻目标网站的服务器负担。使用适当的爬取速率,以防 止对服务器造成过度负担。
- 2. 不滥用爬虫:不要滥用爬虫技术,例如用于非法目的,如网络攻击、诈骗或恶意行为。
- 3. 不侵犯隐私:避免爬取和分发私人或敏感信息,尊重用户的隐私。
- **4. 公开来源**:如果计划使用爬取的数据,确保数据是公开和公开可用的,或者有合法权利使用该数 据。
- **5. 遵循行业标准**: 遵循行业标准和最佳实践,以确保爬虫活动是合法和道德的。

总之,爬虫活动需要谨慎处理,遵守法律规定和道德原则。了解目标网站的政策和法规,采取适当的 措施来保护用户隐私和网站的合法权益,以确保爬虫活动是合法和道德的。

更多 Python 相关干货 内容,扫码领取!!!



# 公众号:涛哥聊Python



# 干货资料领取:

- 1、【优质资料】优质资料合集
- 2、【学习路线】全方位知识点框架
- 3、【问题】Python各领域常见问题
- 4、【面试】面试指南

也欢迎大家围观我的朋友圈,搞搞技术,吹吹牛逼,朋友圈也会发一些外包单,方便自己没时间的时候,小伙伴可以一起利用技术接一些副业项目赚钱!!

添加涛哥 VX: 257735, 围观朋友圈, 一起学 Python



