

10-用程序指挥浏览器

selenium简介

他是一个强大的python库，我们可以通过它控制浏览器，作出自动打开、输入、点击等操作

静态网页和动态网页

1.用html写出的网页，就是静态网页。我们使用BeautifulSoup爬取这类型网页，因为网页源代码中就包含着网页的所有信息，因此，网页地址栏的URL就是网页源代码的URL。

2.要爬取的数据不在HTML源代码中，而是在json中，不能直接使用网址栏的URL，而需要找到json数据的真实URL。这就是一种动态网页。

浏览器总是在向服务器发起各式各样的请求，当这些请求完成后，它们会一起组成开发者工具的Elements中所展示的，渲染完成的网页源代码。

selenium使用场景

selenium的优点：

遇到页面交互复杂或是URL加密逻辑复杂的情况时，我们可以使用selenium打开一个浏览器，等待所有数据都加载到Elements中之后，再把这个网页当做静态网页爬取

selenium的缺点：

由于要真实地运行本地浏览器，打开浏览器以及等待网渲染完成需要一些时间，selenium的工作不可避免地牺牲了速度和更多资源。

selenium的配置

selenium安装

```
1 pip install selenium # Windows电脑安装selenium
2 pip3 install selenium # Mac电脑安装selenium
```

浏览器配置

<http://npm.taobao.org/mirrors/chromedriver/>

本地Chrome浏览器设置方法：

```
1 #从selenium库中调用webdriver模块
2 from selenium import webdriver
3 # 设置引擎为Chrome，从本地打开一个Chrome浏览器
4 driver = webdriver.Chrome()
5 #以上就是浏览器的设置方式：把Chrome浏览器设置为引擎，然后赋值给变量driver。driver
  是实例化的浏览器，在后面的代码中我们都需要添加这些内容。从而使我们能够控制浏览器。
```

selenium获取处理数据用法

获取数据

```
1 # 本地Chrome浏览器设置方法
2 from selenium import webdriver #从selenium库中调用webdriver模块
3 #导入时间模块
4 import time
5 driver = webdriver.Chrome() # 设置引擎为Chrome，真实地打开一个Chrome浏览器
6 driver.get('https://xiaoke.kaikeba.com/example/X-Man/') # 打开网页
7 time.sleep(3)
8 driver.close() # 关闭浏览器
```

get(URL)是webdriver的一个方法，它的使命是为你打开指定URL的网页。

driver是我们实例化后的一个浏览器，我们就是通过driver打开网页，网页中的数据加载到浏览器以后，数据我们就获取到了。

driver.close()是关闭浏览器的方法，我们每次使用webdriver以后，都需要关闭它，不然他会一直占用我们的内存。

解析并提取数据

selenium库也具备解析数据、提取数据的能力，它和BeautifulSoup的底层原理是一致的，不同的是，selenium解析提取的是Elements中的数据，BeautifulSoup解析的是Nextwork中第0个请求的响应

```
1 from selenium import webdriver #从selenium库中调用webdriver模块
2 #导入time模块
3 import time
4 driver = webdriver.Chrome() # 设置引擎为Chrome，从本地打开一个Chrome浏览器
5 driver.get('https://xiaoke.kaikeba.com/example/X-Man/') # 访问页面
6 time.sleep(3) # 等待3秒
7 label = driver.find_element_by_tag_name('label') # 解析网页并提取第一个
  <label>标签
8 print(label.text) # 打印label的文本
9 driver.close() # 关闭浏览器
```

提取数据的方法：

```
1 #通过元素id提取一个或者所有数据
2 driver.find_element_by_id()
3 driver.find_elements_by_id()
```

```
4 #通过元素的class提取一个或者所有数据
5 driver.find_element_by_class_name()
6 driver.find_elements_by_class_name()
7 #通过元素的标签名tag来提取一个或者所有数据
8 driver.find_element_by_tag_name()
9 driver.find_elements_by_tag_name()
10 #通过元素的name提取一个或者所有数据
11 driver.find_element_by_name()
12 driver.find_elements_by_name()
13 #通过链接的部分文本来获取一个或者所有超链接
14 driver.find_element_by_partial_link_text()
15 driver.find_elements_by_partial_link_text()
16 #通过链接的文本来获取一个或者所有超链接
17 driver.find_element_by_link_text()
18 driver.find_elements_by_link_text()
```

解析数据的方法:

```
1 #获取文字
2 webElement.text
3 #带参数(属性名), 可以获取属性值
4 webElement.get_attribute()
```

总结

1.selenium操作数据转换:

- 1.driver.get()获取数据
- 2.driver.find_element_by提取数据
- 3.webelement解析数据

2.selenium配合BeautifulSoup处理数据

```
1 #获取到渲染完整的网页源代码
2 HTML = driver.page_source
3 #page_source获取到的数据已经是str类型, 不需要text转化
4 bs_movies = BeautifulSoup(HTML,'html.parser')
```

3.selenium常用操作

```
1 .send_keys() # 模拟按键输入, 自动填写表单
2 .click() # 点击元素
3 .clear() #清楚元素所有内容
```

实操代码

```
1 # 本地Chrome浏览器设置方法
2 from selenium import webdriver # 从selenium库中调用webdriver模块
3 import time # 调用time模块
4 driver = webdriver.Chrome() # 设置引擎为Chrome, 真实地打开一个Chrome浏览器
5 driver.get('https://xiaoke.kaikeba.com/example/X-Man/') # 访问页面
```

```

6 time.sleep(2) # 暂停两秒, 等待浏览器缓冲
7 teacher = driver.find_element_by_id('teacher') # 找到【请输入你喜欢的老师】下面的输入框位置
8 teacher.send_keys('开课吧') # 输入文字
9 assistant = driver.find_element_by_name('assist') # 找到【请输入你喜欢的助教】下面的输入框位置
10 assistant.send_keys('全都喜欢') # 输入文字
11 time.sleep(2)
12 button = driver.find_element_by_tag_name('button') # 找到【提交】按钮
13 time.sleep(1)
14 #点击提交按钮
15 button.click()
16 time.sleep(5)
17 #关闭浏览器
18 driver.close()

```

qq音乐获取歌曲《说好不哭》全部精彩评论

```

1 # 这段代码在本地执行
2 from selenium import webdriver #从selenium库中调用webdriver模块
3 import time
4 driver = webdriver.Chrome() # 设置引擎为Chrome, 从本地打开一个Chrome浏览器
5 driver.get('https://y.qq.com/n/yyqq/song/001qvvgF38HVc4.html') # 访问页面
6 time.sleep(2)
7 while True:
8     try:
9         #找到记载更多的按钮
10         clickformore =
11         driver.find_element_by_class_name('js_get_more_hot')
12         time.sleep(0.5)
13         #点击加载更多
14         clickformore.click()
15     except:
16         #如果出错, 执行except
17         print("不能再点击加载更多啦")
18         break
19 comments =
20 driver.find_element_by_class_name('js_hot_list').find_elements_by_class_name('js_cmt_li') # 使用class_name找到评论
21 print(len(comments))
22 for i in range(len(comments)): # 循环
23     comment = comments[i].find_element_by_tag_name('p') # 找到评论
24     print('评论'+str(i)+':'+comment.text + '\n-----')
25     # 打印评论
26 driver.close() # 关闭浏览器

```

浏览器静默模式设定

```

1 # 本地Chrome浏览器的静默模式设置:
2 from selenium import webdriver #从selenium库中调用webdriver模块

```

```
3 from selenium.webdriver.chrome.options import Options # 从options模块中调用
Options类
4 chrome_options = Options() # 实例化Option对象
5 chrome_options.add_argument('--headless') # 把Chrome浏览器设置为静默模式
6 driver = webdriver.Chrome(options = chrome_options) # 设置引擎为Chrome，在后
台默默运行
```