

2025年春

程序设计实习：Python程序设计

第十六讲 网络爬虫设计

刘家瑛

liujiaying@pku.edu.cn



01 | 爬虫的用途和原理



爬虫的用途

- 在网络上搜集数据 (例如: 搜索引擎)
- 模拟浏览器快速操作 (抢票, 抢课, 抢挂号...)
- 模拟浏览器操作, 替代填表等重复操作



最基本的爬虫写法

数据获取型爬虫的本质就是自动获取网页并抽取其中的内容

- ① 手工找出合适的URL (网址)
- ② 用浏览器手工查看URL对应的网页, 并查看网页源码, 找出包含想要的内容 (文件名, 链接等) 的字符串的模式
- ③ 程序中获取URL对应的网页
- ④ 程序中用正则表达式或BeautifulSoup库或通过浏览器查找元素功能抽取网页中想要的内容并保存



示例：获取百度图片的搜索结果图片

1. 在百度图片敲关键字“**desk**”进行搜索

2. 搜索后看浏览器地址栏的地址：

<https://image.baidu.com/search/index?tn=baiduimage&ps=1&ct=201326592&lm=-1&cl=2&nc=1&ie=utf-8&dyTabStr=MCwxMiwzLDEsMiwxMyw3LDYsNSw5&word=desk>

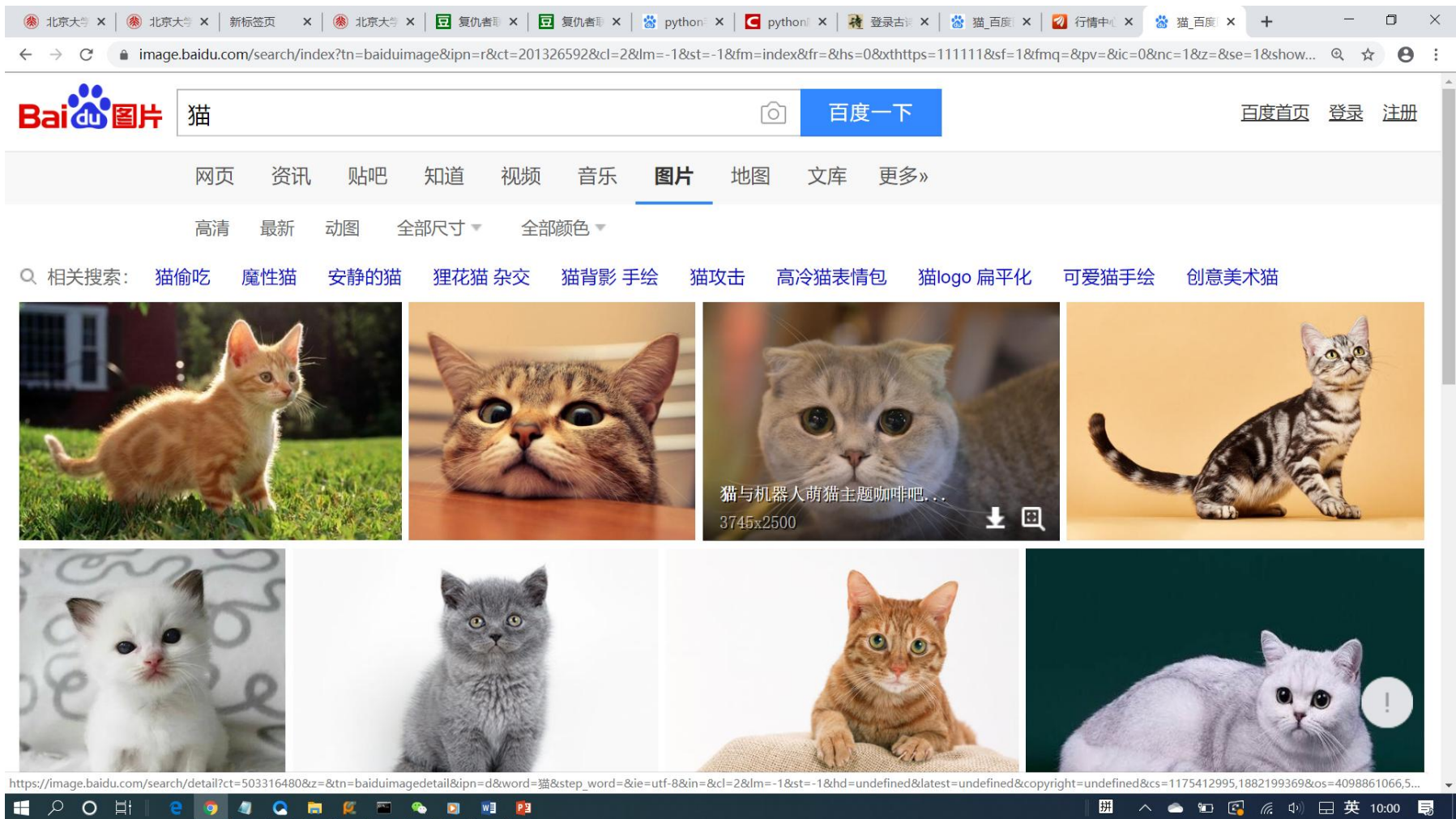


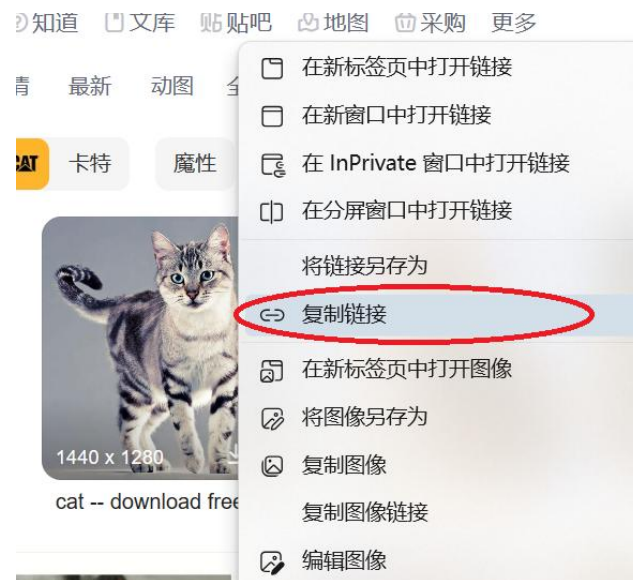
示例：获取百度图片的搜索结果图片

3. 猜测只要在浏览器输入下面地址的红色部分, 替换蓝色部分的单词, 就能搜到图片:

[https://image.baidu.com/search/index?tn=baiduimage&ps=1
&ct=201326592&lm=-1&cl=2&nc=1&ie=utf-
8&dyTabStr=MCwxMiwzLDEsMiwxMyw3LDYsNSw5&word=猫](https://image.baidu.com/search/index?tn=baiduimage&ps=1&ct=201326592&lm=-1&cl=2&nc=1&ie=utf-8&dyTabStr=MCwxMiwzLDEsMiwxMyw3LDYsNSw5&word=猫)







鼠标右键点“下载”图标, 在弹出的
菜单上点“复制链接”复制图片的**URL**

复制出来的**图片**地址:

<https://img2.baidu.com/it/u=519022685,452596173&fm=253&fmt=auto&app=120&f=JPEG?w=800&h=800>

此链接是百度保存的缩略图的网址



示例：获取百度图片的搜索结果图片

4. 在搜索结果界面, 用浏览器查看源码

5. 在源码中查找

`u=519022685,452596173&fm=253&fmt=auto&app=120&f=JPEG?w=800&h=800`

\u0026 是 & 的Unicode转义序列:

- \u 表示Unicode转义

- 0026 是 & 的Unicode码点 (十六进制)

找到:

...

`{"thumburl":"https://img2.baidu.com/it/u=519022685,452596173\u0026fm=253\u0026fmt=auto\u0026app=120\u0026f=JPEG?w=800\u0026h=800",`

...

可以用正则表达式提取图片链接



```
import re                #使用正则表达式
import requests          #request库用于获取网络资源 pip install requests
def getHtml(url):        #获取网址为url的网页, 返回值是个字符串
#具体实现略, 后面再讲
def getBaiduPictures(word, n):
#下载n个百度图片搜来的关于word的图片保存到本地
    url = "https://image.baidu.com/search/index?tn=baiduimage&ipn=r
    &ct=201326592&cl=2&lm=-1&st=1&fm=index&fr=&hs=0
    &xthttps=111111&sf=1&fmq=&pv=&ic=0&nc=1&z=&se=1&showtab=0&fb=0
    &width=&height=&face=0&istype=2&ie=utf-8&word="
    url += word
    html = getHtml(url)
    pt = '\thumburl\":.*?\\"(.*)\\"' #正则表达式, 用于寻找图片网址
    i = 0
#thumburl": "https://img2.baidu.com/it/u=519022685,452596173\u002
6fm=253\u0026fmt=auto\u0026app=120\u0026f=JPEG?w=800\u0026h=800",
```



```
for x in re.findall(pt, html):  
    print(x)  
    x = x.replace(r'\u0026', '&')  
    try:  
        r = requests.get(x, stream=True) #获取x对应的网络资源  
        f = open('{0}{1}.jpg'.format(word, i), "wb")  
        # "wb"表示二进制写方式打开文件  
        f.write(r.content) #图片内容写入文件  
        f.close()  
        i = i + 1  
    except Exception as e :  
        pass  
    if i >= n:  
        break  
  
getBaiduPictures("猫", 5)  
getBaiduPictures("熊猫", 5)
```



02 | 用 requests 库获取网页



□ 用requests.get获取网页

```
import requests
```

```
#requests库用于获取网络资源, pip install requests
```

requests库：Python中用于发送 HTTP 请求的第三方库，
需通过 `pip install requests` 安装

作用：发起 GET 请求并获取网页内容



□ 用requests.get获取网页

```
def getHtml(url):    #获取网址url的网页
    import requests
    fakeHeaders = {'User-Agent':    #用于伪装浏览器发送请求
        'Mozilla/5.0 (Windows NT 10.0; Win64; x64) \
        AppleWebKit/537.36 (KHTML, like Gecko) \
        Chrome/81.0.4044.138 Safari/537.36 Edg/81.0.416.77',
        'Accept': 'text/html, application/xhtml+xml, */*'}
    try:
        r = requests.get(url, headers = fakeHeaders)
        r.encoding = r.apparent_encoding    #确保网页编码正确
        return r.text    #返回值是个字符串, 内含整个网页内容
    except Exception as e:
        print(e)
        return ""

#用法: html = getHtml("http://openjudge.cn")
```



□ Http请求头 fakeHeaders

```
fakeHeaders = { 'User-Agent':    #用于伪装浏览器发送请求  
    'Mozilla/5.0 (Windows NT 10.0; Win64; x64) \\  
    AppleWebKit/537.36 (KHTML, like Gecko) \\  
    Chrome/81.0.4044.138 Safari/537.36 Edg/81.0.416.77',  
    'Accept': 'text/html,application/xhtml+xml,*/*' }
```

User-Agent

作用：告诉服务器你使用的浏览器和操作系统信息

含义：

Mozilla/5.0：历史遗留字段，现代浏览器都会带这个标识

Windows NT 10.0; Win64; x64：表示运行在 Windows 10 64位系统

AppleWebKit/537.36：浏览器引擎 (Chrome/Edge/Safari都使用它)

Chrome/81.0.4044.138：基于Chromium的浏览器 (如Chrome或Edge)

Safari/537.36：兼容 Safari 的渲染引擎

Edg/81.0.416.77：表示这是 Microsoft Edge 浏览器 (版本 81)



□ Http请求头 fakeHeaders

```
fakeHeaders = { 'User-Agent':    #用于伪装浏览器发送请求  
    'Mozilla/5.0 (Windows NT 10.0; Win64; x64) \\  
    AppleWebKit/537.36 (KHTML, like Gecko) \\  
    Chrome/81.0.4044.138 Safari/537.36 Edg/81.0.416.77',  
    'Accept': 'text/html,application/xhtml+xml,*/*' }
```

Agent

作用：告诉服务器客户端可以接收哪些类型的响应数据

含义：

text/html：优先接收 HTML 格式的数据。

application/xhtml+xml：也可以接收 XHTML 格式

/*/*：如果前两种不可用，可以接受任何类型的数据



□ Http请求头 fakeHeaders

```
fakeHeaders = {'User-Agent': #用于伪装浏览器发送请求  
               'Mozilla/5.0 (Windows NT 10.0; Win64; x64) \\  
               AppleWebKit/537.36 (KHTML, like Gecko) \\  
               Chrome/81.0.4044.138 Safari/537.36 Edg/81.0.416.77',  
               'Accept': 'text/html,application/xhtml+xml,*/*\')}
```

Agent

作用：告诉服务器客户端可
含义：

text/html：优先接收 HT

application/xhtml+xml

/：如果前两种不可用，任

为什么需要 fakeHeaders?

1. 避免被反爬虫拦截

- 某些网站会检查**User-Agent**，如果发现是Python的默认UA，可能会拒绝响应或返回错误数据
- 通过伪装成浏览器，可以绕过简单的反爬机制

2. 模拟真实用户访问

- 有些网站会根据**Accept**返回不同的内容（如PC版或移动版页面），设置合理的**Accept**可以确保获取正确的数据

3. 兼容性

- 部分网站依赖**User-Agent**判断浏览器类型，如果没有正确的UA，可能导致JavaScript/CSS加载异常

□ 用requests库获取网页的优势和局限:

优势

- 相比其它方法, 速度快几倍
- 安装简单, 分发容易

局限

- 容易被反爬虫手段屏蔽
- 不能获取包含javascript生成的动态网页



□ 弥补requests不足的其它获取网页的办法

- 使用selenium库

慢, 很容易被反爬, 且已经被许多网站反爬, 网上各种对付反爬的解决办法基本不管用, 不推荐

- 使用puppeteer库

快, 暂未被许多网站反爬, 强烈推荐



03 | 用 selenium 库 获取 网页



□ selenium环境安装

```
pip install selenium
```

- 需要**chrome**浏览器或**firefox**浏览器, 此外还需要下载**chrome**驱动程序 (**chromedriver.exe**) 或**firefox**驱动程序 (**geckodriver.exe**)
- 安装方法:

<https://googlechromelabs.github.io/chrome-for-testing/>

https://blog.csdn.net/m0_49449205/article/details/123971402?spm=1001.2014.3001.5501



□ 用selenium获取网页

```
def getHtml(url):#暂时适用于百度图片搜索,慢且容易被反爬不是推荐的做法,随便看看
    from selenium import webdriver    #需要pip install selenium
    from selenium.webdriver.chrome.options import Options
    from selenium.webdriver.chrome.service import Service
    options = Options()    #浏览器选项
    #等价于 options = webdriver.chrome.options.Options()
    options.add_argument('--headless') #规定chrome浏览器隐身模式运行
    options.add_argument('--disable-gpu')
    #禁止chrome使用gpu加速,能快点
    service = Service(executable_path='./chromedriver.exe')
    driver = webdriver.Chrome( service=service, options=options )
    #driver就是个chrome浏览器,需要下载安装chrome驱动器chromedriver.exe
    driver.get(url)                #浏览器装入网页
    html = driver.page_source      #网页源代码
    driver.close()                 #关闭浏览器
    driver.quit()                  #退出
    return html                    #返回字符串
```



04 | 用 pyppeteer 库 获取 网页



□ pyppeteer 来历

- puppeteer是谷歌公司推出的可以控制Chrome浏览器的一套编程工具
- 一个日本工程师以此为基础推出了Python版本叫pyppeteer
- pyppeteer的官网在: <https://pypi.org/project/pyppeteer/>



□ pyppeteer工作原理

- 启动一个浏览器**Chromium**，用浏览器装入网页；浏览器可以用**无头模式 (headless)**，即隐藏模式启动，也可以显式启动
- 从浏览器可以获取网页源代码，若网页有**javascript**程序，获取到的是**javascript**被浏览器执行后的网页源代码
- 可以向浏览器发送命令，模拟用户在浏览器上键盘输入，鼠标点击等操作，让浏览器转到其它网页
- **selenium**原理及功能和**pyppeteer**一样



□ pyppeteer环境安装

`pip install pyppeteer`

要求Python版本 ≥ 3.6

必须下载并安装特殊版本的谷歌浏览器Chromium

- 可以将Chromium压缩包随便解压在哪个文件夹, 然后在程序指明其中chrome.exe的位置,
- 也可以将Chromium解压到pyppeteer的安装文件夹下面, 这个文件夹通常类似:
C:\Users\username\AppData\Local\pyppeteer\pyppeteer\local-chromium\588429
- 把username要换成自己的windows用户名, 588429这里可能是别的数;
- 将Chromium压缩包里面chrome-win32文件夹整个放在上面那个文件夹里面就行

强烈推荐作为爬虫首选工具!!!



□ 预备知识: 协程 (Coroutine)

协程是一种轻量级的线程, 可以在单个线程内实现并发执行
与线程不同, 协程的切换由程序控制 (而非操作系统调度),
避免了线程切换的开销

特点:

- 非抢占式: 协程主动让出执行权 (通过 `await`), 而不是被强制中断
- 协作式调度: 协程之间需要显式协作切换
- 适合I/O密集型任务: 如网络请求、文件读写等



□ 预备知识: 协程的关键组件

(1) `async` 和 `await` 关键字

`async`: 声明一个函数为异步函数（协程函数）

```
async def fetch_data():  
    return "data"
```

`await`: 挂起当前协程, 等待另一个协程或异步操作完成

```
async def main():  
    data = await fetch_data() # 等待fetch_data()执行完毕  
    print(data)
```

(2) 事件循环 (Event Loop)

协程的运行依赖于事件循环, 负责调度协程的执行



□ 预备知识: 协程的关键组件

(2) 事件循环 (Event Loop)

协程的运行依赖于事件循环, 负责调度协程的执行

```
import asyncio

async def task():
    print("协程开始")
    await asyncio.sleep(1)  # 模拟I/O操作
    print("协程结束")

# 创建事件循环并运行协程
asyncio.run(task())  # Python 3.7+
```



□ 预备知识: 协程

协程就是前面加了 '**async**' 的函数 (从Python 3.6开始有)

```
async def f()  
    return 0
```

- 调用协程时, 必须在函数名前面加 '**await**'

```
await f()
```

- 协程只能在协程里面调用, 即**await**语句只能出现在协程里面
- 协程是一种特殊的函数, 多个协程可以并行

pyppeteer中的所有函数都是协程, 调用时前面都要加 **await**, 且只能在协程中调用

- 初用协程, 经常因为调用XXX时忘了加 **await** 导致下面错误:

Runtime Warning: coroutine 'XXX' was never awaited



□ 用pyppeteer获取网页

def getHtml(url): #暂时适用于百度图片搜索

import asyncio

#Python 3.6之后自带的协程库

import pyppeteer as pp

async def asGetHtml(url):

browser = await pp.launch({

启动Chromium,

page = await browser.newPage()

在浏览器中打开url

await page.goto(url)

Win64; x64) AppleWebKit/537.36

Chrome/78.0.3930.162 Safari/537.36

await page.evaluateOnNewDocument(

'() =>{ Object.defineProperty(navigator, \

{ webdriver:{ get: () => false } }) }')

#反反爬措施

1. navigator.webdriver

- 是浏览器暴露的一个属性, 默认情况下, 当通过 Selenium/Puppeteer/Pyppeteer 等自动化工具控制浏览器时, 该属性会返回 true。
- 许多网站会检测 navigator.webdriver 来判断当前访问是否由自动化脚本发起, 如果是则可能拒绝服务或返回假数据。

2. Object.defineProperty

- 通过重新定义 navigator.webdriver 的 getter 方法, 强制其始终返回 false。
- 这样网站检测时会被欺骗, 认为这是普通用户的手动操作。



□ 用pyppeteer获取网页

```
await page.goto(url)  # 装入url对应的网页
text = await page.content()
# page.content就是网页源代码字符串
await browser.close()  # 关闭浏览器
return text
#速度大约比用requests.get慢5,6倍
```

```
loop = asyncio.new_event_loop() # 创建新事件循环
asyncio.set_event_loop(loop)    # 设置为当前线程的事件循环
html = loop.run_until_complete(asGetHtml(url))
#返回值就是asGetHtml的返回值
return html
```

可以改进程序,只需要启动一次浏览器,生成一个page对象即可,以后获取不同网页都用相同page对象,所有事情完成后才关闭浏览器



□ 用pyppeteer获取网页

launch的其它参数

```
browser = await launch(headless=False,  
    executablePath = "c:/tmp/chrome-win32/chrome.exe",  
    userDataDir = "c:/tmp")
```

- **executablePath**: 如果Chromium没有安装在默认文件夹下面, 则需要指定其位置
- **userDataDir**: **userDataDir**指明用来存放浏览器工作期间存放临时文件的文件夹; 不是必须, 能够防止可能出现的莫名其妙的错误



05 | 用 BeautifulSoup库分析网页



分析并提取网页内容的三种方式

1. 正则表达式 (速度最快, 但适应变化略差)
2. **BeautifulSoup**库 (速度是正则表达式的约几分之一)
3. **selenium**或**pyppeteer**的中的浏览器对象的查找元素函数 (速度是正则表达式的约百分之一, 用在需要模拟在网页中进行输入, 点击按钮等操作的时候)



网页html文档中的tag

```
<div id="siteHeader" class="wrapper">
  <h1 class="logo">
    <a href="http://openjudge.cn/">OpenJudge<span>开放的在线程序评测平台</span></a>
  </h1>
  <div id="topsearch">
    <ul id="userMenu">
      <li class="first"><a href="http://openjudge.cn/">首页</a></li>
      <li><a href="http://openjudge.cn/user/2312/">个人首页</a></li>
      <li><a href="http://openjudge.cn/groups">小组</a></li>
      <li><a href="http://openjudge.cn/settings/">设置</a></li>
      <li><a href="http://openjudge.cn/messages/">信箱(43)</a></li>
      <li><a href="http://openjudge.cn/auth/logout/?1513652042" title="登出">登出</a></li>
    </ul>
  </div>
</div>
```

Diagram illustrating HTML tag components:

- attr** (Attribute): Points to `id="siteHeader"` and `class="wrapper"` in the `<div>` tag.
- text** (Text): Points to the content `OpenJudge` and `开放的在线程序评测平台` within the `<a>` tag.
- attr** (Attribute): Points to `href="http://openjudge.cn/"` and `title="登出"` in the `<a>` tag.



html文档中的tag

□ tag格式通常为(少数没有正文和</x>)：

```
<X attr1='xxx' attr2='yyy' attr3='zzz' ...>
```

n n n n n n n n n n n n

</x>

- **X**: **tag**的名字 (**name**)
- **attr1, attr2, ...**: **tag**的属性 (**attr**)=后面跟着属性的值
- **nnnnnnnnnnnnnnnnnnnn**: **tag**的正文 (**text**)



html文档中的tag

□ 例如：

```
<a href="www.sohu.com" id='mylink'>搜狐网</a>
```

- **a**: tag的名字 (**name**)
- **href, id**: tag的属性 (**attr**), =后面跟着属性的值
- 搜狐网: tag的正文 (**text**)



tag可以嵌套

```
<div id="siteHeader" class="wrapper">  
  <h1 class="logo">  
    <div id="topsearch">  
      <ul id="userMenu">  
        <li ><a href="http://openjudge.cn/">首页</a></li>  
      </div>  
    </div>
```



用 BeautifulSoup库分析html

□ 安装：

```
pip install beautifulsoup4
```

□ 导入：

```
import bs4
```



用 BeautifulSoup库分析html

□ 使用：

- 1) 将html文档装入一个BeautifulSoup对象X
- 2) 用X对象的find, find_all等函数去找想要的tag对象
- 3) 对找到的tag对象，还可以用其find, find_all函数去找它内部包含（嵌套）的tag对象
- 4) tag对象的text就是该对象里的正文(text)，tag对象也可以看作是一个字典，里面包含各种属性(attr)及其值



把html文档载入BeautifulSoup对象

□ 方法1 html文档来自字符串:

```
str = '''
<div id="siteHeader" class="wrapper">
    <h1 class="logo">
    <div id="topsearch">
        <ul id="userMenu">
            <li><a href="http://openjudge.cn/">首页</a></li>
        </div>
    </div>
</div>
'''
```

带href的<a>都是链接, 上面“首页”是链接文字, href后面

http://openjudge.cn是链接地址

```
soup = bs4.BeautifulSoup(str, "html.parser")
```

```
print(soup.find("li").text)      #>> 首页
```



把html文档载入BeautifulSoup对象

□ 方法2 html文档来自于文件:

```
soup = bs4.BeautifulSoup(  
    open("c:\\tmp\\test.html", "r", encoding="utf-8"),  
    "html.parser")
```



把html文档载入BeautifulSoup对象

□ 方法3 html文档来自于给定网址:

```
import requests
def getHtml(url):
    #获得html文本
    try:
        r = requests.get(url)
        r.raise_for_status()
        r.encoding = r.apparent_encoding
        return r.text
    except:
        return ""
```



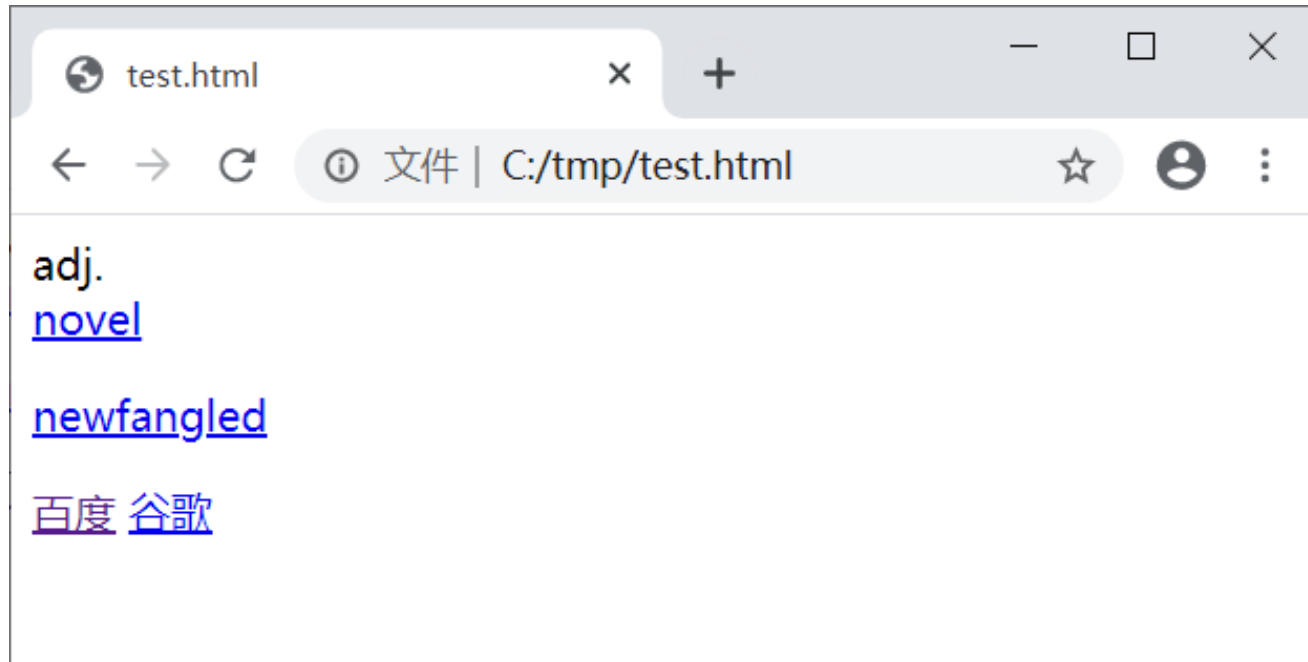
把html文档载入BeautifulSoup对象

□ 方法3 html文档来自于给定网址:

```
html = getHtml("https://cn.bing.com/dict/search?q=new")  
soup = bs4.BeautifulSoup(html, "html.parser")
```



c:\tmp\test.html:



用 BeautifulSoup对象寻找想要的tag

c:\tmp\test.html:

```
<!DOCTYPE HTML>
<html>
<body>
<div id="synoid" style="display:block;">
  <div class="df_div2">
    <div class="de_title1">adj.</div>
    <div class="col_fl">
      <a h="ID=Dictionary,5237.1"
        href="https://cn.bing.com/dict/search?q=novel">
        <span class="p1-4">novel</span>
      </a>
    <p>
```



用 BeautifulSoup对象寻找想要的tag

c:\tmp\test.html:

```
<a h="ID=Dictionary,5238.1"
href="https://cn.bing.com/dict/search?q=newfangled">
<span class="p1-4">newfangled</span>
</a>
```

```
</div>
```

```
<a href="http://www.baidu.com" id="searchlink1" class="sh1">百度</a>
<a href="http://www.google.com" id="searchlink1" class="sh2">谷歌</a>
</div>
```

```
</div>
```

```
</body>
```

```
</html>
```



用 BeautifulSoup 对象寻找想要的 tag

```
import bs4
soup = bs4.BeautifulSoup(open("c:\\tmp\\test.html",
                              encoding = "utf-8"), "html.parser")


diva = soup.find("div", attrs={"id": "synoid"})



#寻找名为"div", 且具有 值为"synoid"的属性"id"的tag



if diva != None: #如果找到



```
 for x in diva.find_all("span", attrs={"class": "p1-4"}):
 print(x.text) #在diva内部继续找
```



```
 for x in diva.find_all("a", attrs={"id": "searchlink1"}):
 print(x.text)
```



```
 x = diva.find("a", attrs={"id": "searchlink1", "class": "sh2"})
```



```
 if x != None:
```



```
 print(x.text)
```



```
 print(x["href"])
```



```
 print(x["id"])
```



50



A collage of three images: a traditional Chinese building with a tiled roof, a modern building with a glass facade, and a classical building with columns.


```


用 BeautifulSoup对象寻找想要的tag

输出：

novel

newfangled

百度

谷歌

谷歌

<http://www.google.com>

searchlink1



用 BeautifulSoup对象寻找想要的tag

如果别处没有和 `<div id="synoid" style="display:block;">`
内部的模式相似的东西,也可以不用先找这个"synoid" tag

```
import bs4
soup = bs4.BeautifulSoup(
    open("c:\\tmp\\test.html", encoding="utf-8"), "html.parser")
for x in soup.find_all("span", attrs={"class": "p1-4"}):
    print(x.text)
for x in soup.find_all("a", attrs={"id": "searchlink1"}):
    print(x.text)
x = soup.find("a", attrs={"id": "searchlink1", "class": "sh2"})
if x != None:
    print(x.text)
    print(x["href"])
    print(x["id"])
```



06 | 实例：爬取每日股票交易信息



爬取每日股票交易信息

浏览器地址栏: 深振业A(000006)股票 × 上证所有股票代码及名 × view-source:https://w × 创业板所有股票代码及 × 创业板所有股票代码及 ×

URL: /gupiao/list_sh.html

导航栏: 首页 | 概念股大全 | 个股资讯 | A股要闻 | 龙虎榜 | 新股申购 | 新股中签查询 | 公司大全 | 近期热点 | 大盘行

您的位置: 首页 » 股票大全 » 上证股票代码一览表

搜索框: 搜搜看 有没有你想找的

上证所有股票代码及名称

包钢股份(600010)	四川路桥(600039)	保利地产(600048)	九鼎投资(600053)
国投资本(600061)	皖维高新(600063)	凤凰光学(600071)	康欣新材(600076)
宋都股份(600077)	澄星股份(600078)	人福医药(600079)	东风科技(600081)
海泰发展(600082)	同济堂(600090)	恒力石化(600346)	明星电力(600101)
北方稀土(600111)	浙江东日(600113)	东方航空(600115)	浙江东方(600120)
商赢环球(600146)	航天机电(600151)	华升股份(600156)	巨化股份(600160)
天坛生物(600161)	香江控股(600162)	太原重工(600169)	美都能源(600175)
伊力特(600197)	金种子酒(600199)	江苏吴中(600200)	中金公司(601995)
福日电子(600203)	有研新材(600206)	安彩高科(600207)	新湖中宝(600208)
亚星客车(600213)	全柴动力(600218)	海航控股(600221)	鲁商发展(600223)



□ 爬取每日股票交易信息

东方财富网每日股票交易信息单只股票：

`quote.eastmoney.com/sh600000.html` 上海证券

`quote.eastmoney.com/sz000017.html` 深圳或创业板



□ 用pypetter爬取每日上证股票交易信息

单只股票：`quote.eastmoney.com/sh600000.html`



该页面查看源代码, 看不到 12.17, 12.51 等交易数据; 说明源代码里面包含javascript程序, 浏览器执行javascript程序以后, 才能得到显示的页面; 因此python程序需要在取到网页后, 还要执行里面的javascript程序, 才能得到股票数据

用requests.get无法得到显示的网页, 必须用selenium或者pypeteer



用 pyppeteer+正则表达式 爬取每日上证股票交易信息

```
html = getHtml("https://quote.eastmoney.com/sh600000.html")
```

```
# pypeteer版
```

```
# print(html) 先打印出来, 将打印结果拷贝粘贴到记事本,
```

```
# 查找关心的数据如12.17在哪里
```

```
class="quote3l_c"><div class="brief_info_c"><table><tbody><tr><td class="n">今开: </td><td><span><span  
class="price_up blinkred">12.17</span></span></td><td class="n">最高: </td><td><span><span  
class="price_up blinkred">12.51</span></span></td><td class="n">涨停: </td><td><span><span  
class="price_up blinkred">13.34</span></span></td><td class="n">换手: </td><td><span><span
```

正则表达式匹配的格式:

```
<td class="n">XXXX</td><td><span><span class="XXXX">XXXX</span></span></td></span></td>
```



用 pyppeteer+正则表达式 爬取每日上证股票交易信息

```
html = getHtml("https://quote.eastmoney.com/sh600000.html")

# pypeteer版

# print(html) 先打印出来, 将打印结果拷贝粘贴到记事本,
# 查找关心的数据如12.17在哪里

pt = r'<td class="n">([^<]*)</td><td><span><span class=[^<]*>
([^<]*)</span></span></td>'

for x in re.findall(pt, html, re.DOTALL):
    if (x[1] != ""):
        print(x[0], x[1])
```

re.DOTALL: 常规情况下.匹配任意字符但不包括换行符, 使用re.DOTALL后, .可以匹配换行符

正则表达式pt对应找到的关心的数据所在字符串的模式:

```
<td class="n">今开:</td><td><span><span class="price_up blink
red">12.17</span></span></td>
```



用 pyppeteer+正则表达式 爬取每日上证股票交易信息

输出结果:

今开: 6.61

最高: 6.63

涨停: 7.27

换手: 0.05%

成交量: 14.81万

总市值: 1934亿

昨收: 6.61

最低: 6.53

跌停: 5.95

量比: 1.86

成交额: 9724万

市净: 0.32

流通市值: 1934亿



用pyppeteer+BeautifulSoup爬取每日上证股票交易信息

网页内容:

```
<div class="brief_info_c">
  <table>
    <tbody>
      <tr><td class="n">今开: </td><td><span><span class="price_up blinkred">11.73</span></span></td>
      <td class="n">最高: </td><td><span><span class="price_up blinkred">12.15</span></span></td>
      <td class="n">涨停: </td><td><span><span class="price_up blinkred">12.88</span></span></td>
      <td class="n">换手: </td><td><span><span class="price_draw blinkblue">0.29%</span></span></td>
      <td class="n">成交量: </td><td><span><span class="price_draw blinkblue">85.04万</span></span></td>
      <td class="n">市盈(动)
        <span class="title_help">
          <div class="ssy">
            <ul>
              <li class="ssy_t">·动态市盈率: <span class="price_draw">5.04</span></li>
              <li>总市值除以全年预估净利润, 例如当前一季度净利润1000万, 则预估全年净利润4000万</li>
              <li class="ssy_t">·静态市盈率: <span class="price_draw">7.83</span></li>
              <li>总市值除以上一年度净利润</li>
              <li class="ssy_t">·滚动市盈率: <span class="price_draw">7.80</span></li>
              <li>最新价除以最近4个季度的每股收益</li>
            </ul>
          </div>
        </span></td><td><span><span class="price_draw blinkblue">5.04</span></span></td>
      <td class="n">总市值: </td><td><span><span class="price_draw blinkblue">3546亿</span></span></td>
    </tr>
```



用pyppeteer+BeautifulSoup爬取每日上证股票交易信息

```
html = getHtml("https://quote.eastmoney.com/sh600000.html")
#要用pyppeteer版的getHtml
#print(html) #下一步编程前可以先打印出来,
#将打印结果拷贝粘贴到记事本, 查找关心的数据如12.17在哪里
soup = bs4.BeautifulSoup(html, "html.parser")
diva = soup.find("div", attrs={"class": "brief_info_c"})
titles = diva.find_all("td", attrs = {"class": "n"})
for t in titles:
    sp = t.findNext("span") #从当前标题<td>向后查找第一个<span>
    if 'class' not in sp.attrs:
        print(t.text, sp.text)
#print(sp) 可以看到sp的整个tag的描述, 如:
<span><span class="price_up blinkred">0.01</span>
```



用pyppeteer+浏览器查找元素功能爬取每日上证股票交易信息

https://quote.eastmoney.com/sh600000.html

东方财富网

东方财富

天天基金网

东方财富证券

东方财富期货

Choice数据

股吧

登录

新股申购

新股日历

资金流向

AH股比价

主力排名

板块资金

个股研报

行业研报

盈利预测

千股千评

年报季报

龙虎榜单

限售

沪股通

有额度

资金流入-16.27亿

深股通

有额度

资金流入-29.26亿

港股通

有额度

资金流入-16.27亿

上证指数

600000

更名

两融

融资融券

龙虎榜单

个股研报

行业研报

盈利预测

千股千评

年报季报

龙虎榜单

限售

6.62↑

0.01

0.15%

今开: 6.61

昨收: 6.60

6.61

6.62

6.63

6.64

6.65

6.66

6.67

6.68

6.69

6.70

6.71

6.72

6.73

6.74

6.75

6.76

6.77

6.78

6.79

6.80

6.81

6.82

6.83

6.84

6.85

6.86

6.87

6.88

6.89

6.90

6.91

6.92

6.93

6.94

6.95

6.96

6.97

6.98

6.99

7.00

7.01

7.02

7.03

7.04

7.05

7.06

7.07

7.08

7.09

7.10

7.11

7.12

7.13

7.14

7.15

7.16

7.17

7.18

7.19

7.20

7.21

7.22

7.23

7.24

7.25

7.26

7.27

7.28

7.29

7.30

7.31

7.32

7.33

7.34

7.35

7.36

7.37

7.38

7.39

7.40

7.41

7.42

7.43

7.44

7.45

7.46

7.47

7.48

7.49

7.50

7.51

7.52

7.53

7.54

7.55

7.56

7.57

7.58

7.59

7.60

7.61

7.62

7.63

7.64

7.65

7.66

7.67

7.68

7.69

7.70

7.71

7.72

7.73

7.74

7.75

7.76

7.77

7.78

7.79

7.80

7.81

7.82

7.83

7.84

7.85

7.86

7.87

7.88

7.89

7.90

7.91

7.92

7.93

7.94

7.95

7.96

7.97

7.98

7.99

8.00

8.01

8.02

8.03

8.04

8.05

8.06

8.07

8.08

8.09

8.10

8.11

8.12

8.13

8.14

8.15

8.16

8.17

8.18

8.19

8.20

8.21

8.22

8.23

8.24

8.25

8.26

8.27

8.28

8.29

8.30

8.31

8.32

8.33

8.34

8.35

8.36

8.37

8.38

8.39

8.40

8.41

8.42

8.43

8.44

8.45

8.46

8.47

8.48

8.49

8.50

8.51

8.52

8.53

8.54

8.55

8.56

8.57

8.58

8.59

8.60

8.61

8.62

8.63

8.64

8.65

8.66

8.67

8.68

8.69

8.70

8.71

8.72

8.73

8.74

8.75

8.76

8.77

8.78

8.79

8.80

8.81

8.82

8.83

8.84

8.85

8.86

8.87

8.88

8.89

8.90

8.91

8.92

8.93

8.94

8.95

8.96

8.97

8.98

8.99

9.00

9.01

9.02

9.03

9.04

9.05

9.06

9.07

9.08

9.09

9.10

9.11

9.12

9.13

9.14

9.15

9.16

9.17

9.18

9.19

9.20

9.21

9.22

9.23

9.24

9.25

9.26

9.27

9.28

9.29

9.30

9.31

9.32

9.33

9.34

9.35

9.36

9.37

9.38

9.39

9.40

9.41

9.42

9.43

9.44

9.45

9.46

9.47

9.48

9.49

9.50

9.51

9.52

9.53

9.54

9.55

9.56

9.57

9.58

9.59

9.60

9.61

9.62

9.63

9.64

9.65

9.66

9.67

9.68

9.69

9.70

9.71

9.72

9.73

9.74

9.75

9.76

9.77

9.78

9.79

9.80

9.81

9.82

9.83

9.84

9.85

9.86

9.87

9.88

9.89

9.90

9.91

9.92

9.93

9.94

9.95

9.96

9.97

9.98

9.99

10.00

10.01

10.02

10.03

10.04

10.05

10.06

10.07

10.08

10.09

10.10

10.11

10.12

10.13

10.14

10.15

10.16

10.17

10.18

10.19

10.20

10.21

10.22

10.23

10.24

10.25

10.26

10.27

10.28

10.29

10.30

10.31

10.32

10.33

10.34

10.35

10.36

10.37

10.38

10.39

10.40

10.41

10.42

10.43

10.44

10.45

10.46

10.47

10.48

10.49

10.50

10.51

10.52

10.53

10.54

10.55

10.56

10.57

10.58

10.59

10.60

10.61

10.62

10.63

10.64

10.65

10.66

10.67

10.68

10.69

10.70

10.71

10.72

10.73

10.74

10.75

10.76

10.77

10.78

10.79

10.80

10.81

10.82

10.83

10.84

10.85

10.86

10.87

10.88

10.89

10.90

10.91

10.92

10.93

10.94

10.95

10.96

10.97

10.98

10.99

11.00

11.01

11.02

11.03

11.04

11.05

11.06

11.07

11.08

11.09

11.10

11.11

11.12

11.13

11.14

11.15

11.16

11.17

11.18

11.19

11.20

11.21

11.22

11.23

11.24

11.25

11.26

11.27

11.28

11.29

11.30

11.31

11.32

11.33

11.34

11.35

11.36

11.37

11.38

11.39

11.40

11.41

11.42

11.43

11.44

11.45

11.46

11.47

11.48

11.49

11.50

11.51

11.52

11.53

11.54

11.55

11.56

11.57

11.58

11.59

11.60

11.61

11.62

11.63

11.64

11.65

11.66

11.67

11.68

11.69

11.70

11.71

11.72

11.73

11.74

11.75

11.76

11.77

11.78

11.79

11.80

11.81

11.82

11.83

11.84

11.85

11.86

11.87

11.88

11.89

11.90

11.91

11.92

11.93

11.94

11.95

11.96

11.97

11.98

11.99

12.00

12.01

12.02

12.03

12.04

12.05

12.06

12.07

12.08

12.09

12.10

12.11

12.12

12.13

12.14

12.15

12.16

12.17

12.18

12.19

12.20

12.21

12.22

12.23

12.24

12.25

12.26

12.27

12.28

12.29

12.30

12.31

12.32

12.33

12.34

12.35

12.36

12.37

12.38

12.39

12.40

12.41

12.42

12.43

12.44

12.45

12.46

12.47

12.48

12.49

12.50

12.51

12.52

12.53

12.54

12.55

12.56

12.57

12.58

12.59

12.60

12.61

12.62

12.63

12.64

12.65

12.66

12.67

12.68

12.69

12.70

12.71

12.72

12.73

12.74

12.75

12.76

12.77

12.78

12.79

12.80

12.81

12.82

12.83

12.84

12.85

12.86

12.87

12.88

12.89

12.90

12.91

12.92

12.93

12.94

12.95

12.96

12.97

12.98

12.99

13.00

13.01

13.02

13.03

13.04

13.05

13.06

13.07

13.08

13.09

13.10

13.11

13.12

13.13

13.14

13.15

13.16

13.17

13.18

13.19

13.20

13.21

13.22

13.23

13.24

13.25

13.26

13.27

13.28

13.29

13.30

13.31

13.32

13.33

13.34

13.35

13.36

13.37

13.38

13.39

13.40

13.41

13.42

13.43

13.44

13.45

13.46

13.47

13.48

13.49

13.50

13.51

13.52

13.53

13.54

13.55

13.56

13.57

13.58

13.59

13.60

13.61

13.62

13.63

13.64

13.65

13.66

13.67

13.68

13.69

13.70

13.71

13.72

13.73

13.74

13.75

13.76

13.77

13.78

13.79

13.80

13.81

13.82

13.83

13.84

13.85

13.86

13.87

13.88

13.89

13.90

13.91

13.92

13.93

13.94

13.95

13.96

13.97

13.98

13.99

14.00

14.01

14.02

14.03

14.04

14.05

14.06

14.07

14.08

14.09

14.10

14.11

14.12

14.13

14.14

14.15

14.16

14.17

14.18

14.19

14.20

14.21

14.22

14.23

14.24

14.25

14.26

14.27

14.28

14.29

14.30

14.31

14.32

14.33

14.34

14.35

14.36

14.37

14.38

14.39

14.40

14.41

14.42

14.43

14.44

14.45

14.46

14.47

14.48

14.49

14.50

14.51

14.52

14.53

14.54

14.55

14.56

14.57

14.58

14.59

14.60

14.61

14.62

14.63

14.64

14.65

14.66

14.67

14.68

14.69

14.70

14.71

14.72

14.73

14.74

14.75

14.76

14.77

14.78

14.79

14.80

14.81

14.82

14.83

14.84

14.85

14.86

14.87

14.88

14.89

14.90

14.91

14.92

14.93

14.94

14.95

14.96

14.97

14.98

14.99

15.00

15.01

15.02

15.03

15.04

15.05

15.06

15.07

15.08

15.09

15.10

15.11

15.12

15.13

15.14

15.15

15.16

15.17

15.18

15.19

15.20

15.21

15.22

15.23

15.24

15.25

15.26

15.27

15.28

15.29

15.30

15.31

15.32

15.33

15.34

15.35

15.36

15.37

15.38

15.39

15.40

15.41

15.42

15.43

15.44

15.45

15.46

15.47

15.48

15.49

15.50

15.51

15.52

15.53

15.54

15.55

15.56

15.57

15.58

15.59

15.60

15.61

15.62

15.63

15.64

15.65

15.66

15.67

15.68

15.69

15.70

15.71

15.72

15.73

15.74

15.75

15.76

15.77

15.78

15.79

15.80

15.81

15.82

15.83

15.84

15.85

15.86

15.87

15.88

15.89

15.90

15.91

15.92

15.93

15.94

15.95

15.96

15.97

15.98

15.99

16.00

16.01

16.02

16.03

16.04

16.05

16.06

16.07

16.08

16.09

16.10

16.11

16.12

16.13

16.14

16.15

16.16

16.17

16.18

16.19

16.20

16.21

16.22

16.23

16.24

16.25

16.26

16.27

16.28

16.29

16.30

16.31

16.32

16.33

16.34

16.35

16.36

16.37

16.38

16.39

16.40

16.41

16.42

16.43

16.44

16.45

16.46

16.47

16.48

16.49

16.50

16.51

16.52

16.53

16.54

16.55

16.56

16.57

16.58

16.59

16.60

16.61

16.62

16.63

16.64

16.65

16.66

16.67

16.68

16.69

16.70

16.71

16.72

16.73

16.74

16.75

16.76

16.77

16.78

16.79

16.80

16.81

16.82

16.83

16.84

16.85

16.86

16.87

16.88

16.89

16.90

16.91

16.92

16.93

16.94

16.95

16.96

16.97

16.98

16.99

17.00

17.01

17.02

17.03

17.04

17.05

17.06

17.07

17.08

17.09

17.10

17.11

17.12

17.13

17.14

17.15

17.16

17.17

17.18

17.19

17.20

17.21

17.22

17.23

17.24

17.25

17.26

17.27

17.28

17.29

17.30

17.31

17.32

17.33

17.34

17.35

17.36

17.37

17.38

17.39

17.40

17.41

17.42

17.43

17.44

17.45

17.46

17.47

17.48

17.49

17.50

17.51

17.52

17.53

17.54

17.55

17.56

17.57

17.58

17.59

17.60

17.61

17.62

17.63

17.64

17.65

17.66

17.67

17.68

17.69

17.70

17.71

17.72

17.73

17.74

17.75

17.76

17.77

17.78

17.79

17.80

17.81

17.82

17.83

17.84

17.85

17.86

17.87

17.88

17.89

17.90

17.91

17.92

17.93

17.94

17.95

17.96

17.97

17.98

17.99

18.00

18.01

18.02

18.03

18.04

18.05

18.06

18.07

18.08

18.09

18.10

18.11

18.12

18.13

18.14

18.15

18.16

18.17

18.18

18.19

18.20

18.21

18.22

18.23

18.24

18.25

18.26

18.27

18.28

18.29

18.30

18.31

18.32

18.33

18.34

18.35

18.36

18.37

18.38

18.39

18.40

18.41

18.42

18.43

18.44

18.45

18.46

18.47

18.48

18.49

18.50

18.51

18.52

18.53

18.54

18.55

18.56

18.57

18.58

18.59

18.60

18.61

18.62

18.63

18.64

18.65

18.66

18.67

18.68

18.69

18.70

18.71

18.72

18.73

18.74

18.75

18.76

18.77

18.78

18.79

18.80

18.81

18.82

18.8

日历 资金流向 AH股比价 主力排名 板块资金

度 资金流入-16.27亿 | 深股通 有额度

浦发银行

更名 两融 沪主板 沪股通

span.price_draw.blinkblue

今开: 6.61 最高: 6.63 涨停: 7.27

昨收: 6.61 最低: 6.53 跌停: 5.95

F10档案: 操盘必读 公司概况 经营分析 核心题材 股东研究 公司大事 股本结构 财务分析 分红融资

特色数据: 资金流向 主力控盘 机构散户 龙虎榜单 深度数据 大宗交易 智能点评 融资融券 最新业绩

看行情就用Level-2! 5档盘口变千档, 助您轻松看穿主

浦发银行吧 机构散户 精准买卖点 大单成交 盈利预测 机构调研 问董秘 刷新

个股日历 股权质押: 截止2023年11月10日质押总比例0.11%, 质押总股数3213.51万股, 质押总...

盘前 一天 二天 三天 四天 五天 订阅股价提醒 图片版 动图版 极速版

浦发银行[600000] 2023-12-07 11:30 价格: 6.62 涨幅: 0.15% 成交量: 2059



鼠标右键点击“6.61”的tag 在弹出菜单点“复制selector”, 得到该tag的selector标识:

```
#app > div > div > div.zsquote3l.zs_brief  
> div.quote3l_c > div > table > tbody >  
tr:nth-child(1) > td:nth-child(2) > span  
> span
```

添加属性
编辑属性
编辑为 HTML
重复元素
删除元素
剪切

复制
复制元素
复制 outerHTML
复制selector
复制 JS 路径
复制样式
复制 XPath
复制完整的 XPath

粘贴
隐藏元素
强制状态
中断于
以递归方式展开
折叠子项
捕获节点屏幕截图
滚动到视图
焦点
徽章设置...
存储为全局变量

```
<table>  
  <tbody>  
    <tr>  
      <td class="n">  
        <td>  
          <span>  
            <span class="blinkblue">  
              </span>  
            </td>  
            <td class="n">  
              <td>...</td>  
            <td class="n">
```

今开: 6.61 最高: 6.63 涨停: 7.27 换手:

昨收: 6.61 最低: 6.53 跌停: 5.95 量比:

F10档案: 操盘必读 公司概况 经营分析 核

通过pypeteer的Page对象的
querySelector函数可以由selector
(字符串)得到对应的元素

今开 的selector:

```
#app > div > div > div.zsquote3l.zs_brief > div.quote3l_c > div >  
table > tbody > tr:nth-child(1) > td:nth-child(1)
```

6.61 的selector::

```
#app > div > div > div.zsquote3l.zs_brief > div.quote3l_c > div >  
table > tbody > tr:nth-child(1) > td:nth-child(2) > span > span
```

最高 的selector:

```
#app > div > div > div.zsquote3l.zs_brief > div.quote3l_c > div >  
table > tbody > tr:nth-child(1) > td:nth-child(3)
```

6.63 的selector:

```
#app > div > div > div.zsquote3l.zs_brief > div.quote3l_c > div >  
table > tbody > tr:nth-child(1) > td:nth-child(4) > span > span
```



今开: 6.61 最高: 6.63 涨停: 7.27 换手:

昨收: 6.61 最低: 6.53 跌停: 5.95 量比:

F10档案: 操盘必读 公司概况 经营分析 核

通过pypeteer的Page对象的
querySelector函数可以由selector
(字符串)得到对应的元素

昨收 的selector:

```
#app > div > div > div.zsquote3l.zs_brief > div.quote3l_c > div >  
table > tbody > tr:nth-child(2) > td:nth-child(1)
```

第二行的6.61 的selector::

```
#app > div > div > div.zsquote3l.zs_brief > div.quote3l_c > div >  
table > tbody > tr:nth-child(2) > td:nth-child(2) > span > span
```



用pyppeteer+浏览器查找元素功能爬取每日上证股票交易信息

```
import asyncio    #Python 3.6之后自带的协程库
import pyppeteer as pyp
import re

async def antiAntiCrawler(page):    #为page添加反反爬虫手段
    await page.setUserAgent('Mozilla/5.0 (Windows NT 6.1; Win64; x64) \
                             AppleWebKit/537.36 (KHTML, like Gecko) \
                             Chrome/78.0.3904.70 Safari/537.36')
    await page.evaluateOnNewDocument(
        '() =>{ Object.defineProperty(navigator, \
        { webdriver:{ get: () => false } }) }')
```



用pyppeteer+浏览器查找元素功能爬取每日上证股票交易信息

```
async def asGetStockInfo(url):
    # url是"https://quote.eastmoney.com/sh600000.html"
    browser = await pyp.launch(headless=False)
    page = await browser.newPage()
    await antiAntiCrawler(page)
    await page.goto(url)
    html = await page.content()
    for k in range(1,3):
        for i in range(1,14,2):
            selector1 = '#app > div > div > div.zsquote3l.zs_brief ' + \
                '> div.quote3l_c > div > table ' + \
                '> tbody > tr:nth-child(%d) > td:nth-child(%d)' % (k, i)
            selector2 = ('#app > div > div > div.zsquote3l.zs_brief' + \
                '> div.quote3l_c > div > table > tbody > tr:nth-child(%d)' + \
                '> td:nth-child(%d) > span > span') % (k, i+1)
```



用pyppeteer+浏览器查找元素功能爬取每日上证股票交易信息

```
title = await page.querySelector(selector1)
if title is None:
    break
else:
    value = await page.querySelector(selector2)
    obj = await title.getProperty("innerText")
    text = await obj.jsonValue() # 固定写法
    print(text, end = " ")
    obj = await value.getProperty("innerText")
    text = await obj.jsonValue() # 固定写法
    print(text)
await browser.close()
```

#以下是调用协程的固定写法

```
loop = asyncio.new_event_loop()
asyncio.set_event_loop(loop)
url = "https://quote.eastmoney.com/sh600000.html"
loop.run_until_complete(asGetStockInfo(url))
```



07 | 需要登录的爬虫



□ 需要登录的爬虫

- 许多网站需要登录后才能访问其内容

京东、淘宝需要登录才能访问交易记录

openjudge.cn 需要登录才能看提交过的源代码

- 登录操作, 无法用一个url表示出来
- 解决办法之一: 用浏览器模拟登录过程, 输入用户名密码、点登录按钮, 或者程序启动浏览器, 等待手工登录后, 程序再继续爬虫操作 (对有验证码的情况)



爬取Openjudge自己提交通过的所有程序源码



OpenJudge 开放的在线程序评测平台

账号:

密码:

[忘记密码?](#)

[还没有OpenJudge账号?点此注册](#)



OpenJudge 开放的在线程序评测平台

[首页](#) | [个人首页](#) | [小组](#) | [设置](#)

你还没有头像, 设置头像后, 别人会更容易记住你。

正在进行的比赛(90)... (查看全部)

开始时间: 2020-09-17 08:30:00
结束时间: 2020-10-11 23:59:59

 计科19算法分析与设计第1次实验 (焦晓军班)

2	52	282
题目	参与者	提交数

活跃的小组... (查看全部)

 百练

 程序设计实习

 数据结构与算法MOOC

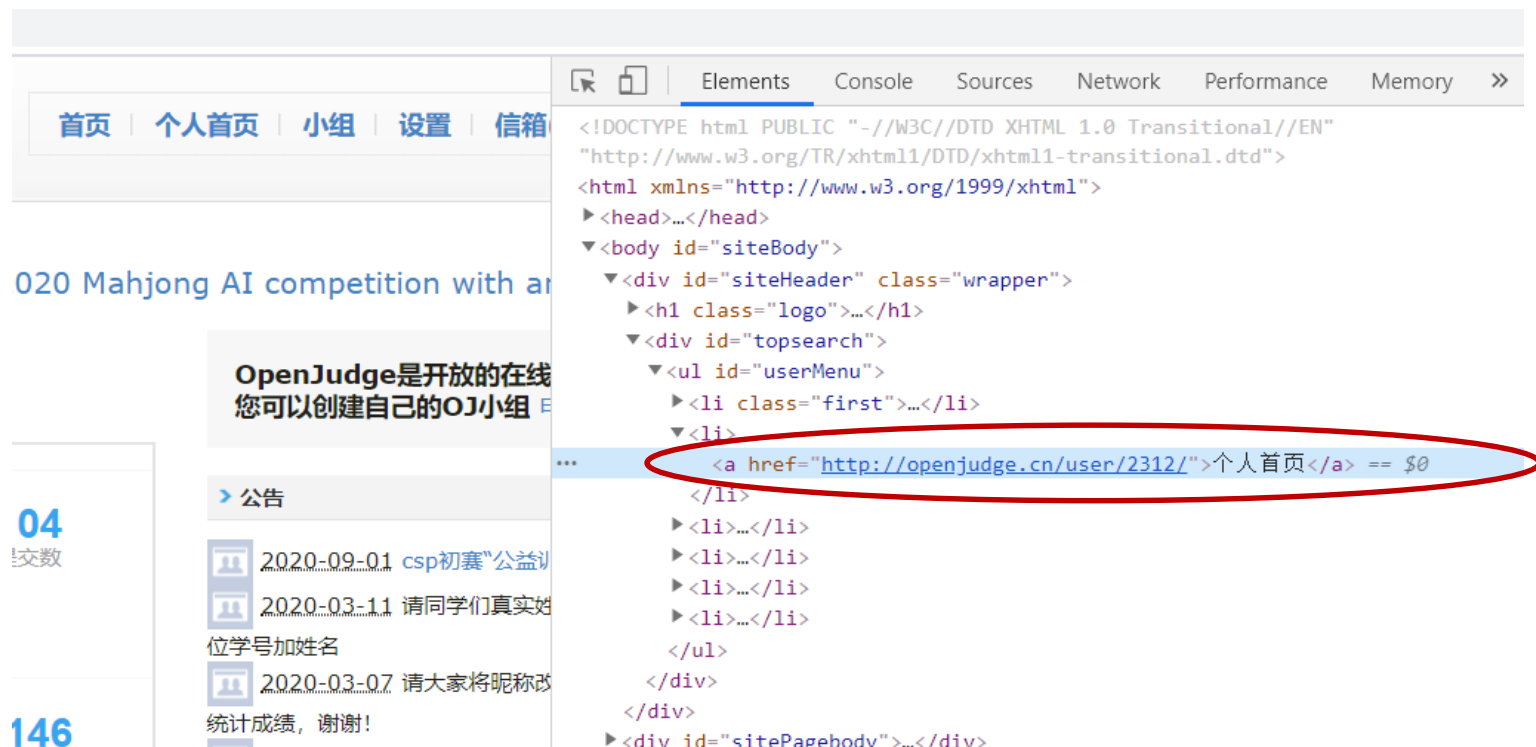
- 程序命令浏览器模拟登录过程, 即输入用户名密码, 点登录按钮
- 或 程序启动浏览器, 等待手工登录后, 程序再继续爬虫操作 (对有验证码的情况, 或者懒得写代码的情况)

- 更高级做法: 不用浏览器, 经数据包分析后, 用requests库进行数据传输进行登录



爬取Openjudge自己提交通过的所有程序源码

鼠标右键点击右上角的“个人首页”，在弹出的菜单上选“检查”：



爬取Openjudge自己提交通过的所有程序源码

点击“个人首页”，进入：

OpenJudge

开放的在线程序评测平台

首页 | 个人首页 | 小组 | 设置 | 信箱(67) | 登出

我加入的小组

最近的提交

百练

北京大学ACM-ICPC竞赛

程序设计实习

PKU Online Judge

程序设计实习实验班

程序设计实习MOOC

比赛	题目	结果	内存	时间	代码长度	语言	提交时间
程序设计 & 算法	018: 整数序列的元素最大跨度值	Accepted	200kB	14ms	281 B	G++	5天前
程序设计 & 算法	018: 整数序列的元素最大跨度值	Accepted	200kB	15ms	281 B	G++	5天前
实用Python程序设计	014: 三角形判断	Accepted	200kB	9ms	190 B	G++	5天前

`Accepted`



爬取Openjudge自己提交通过的所有程序源码

点击某个题的 Accepted 链接, 进入:

题目 排名 状态 统计 提问

#25212869提交状态

状态: Accepted

下载通过码文件

源代码

```
#include<iostream>
#include <stdio>
using namespace std;
int main()
{
    int n,minV , maxV;
    cin >> n >> minV;
    maxV = minV;
```

```
<pre class="sh_python">
n = int(input())
lst = []
for i in range(n):
    s = input().split()
    lst.append((s[0], int(s[1])))
lst.sort(key= lambda x:(-x[1], x[0]))
for x in lst:
    print(x[0], x[1])</pre>
```



❑ pyppeteer 爬取 Openjudge 自己提交通过的所有程序源码

```
import asyncio
import pyppeteer as pyp

async def antiAntiCrawler(page):
    # 为page添加反反爬虫手段
    await page.setUserAgent('Mozilla/5.0 (Windows NT 6.1; Win64; x64) \
                             AppleWebKit/537.36 (KHTML, like Gecko) \
                             Chrome/78.0.3904.70 Safari/537.36')
    await page.evaluateOnNewDocument(
        '() =>{ Object.defineProperty(navigator, \
        { webdriver:{ get: () => false } }) }')
```



❏ pyppeteer 爬取 Openjudge 自己提交通过的所有程序源码

```
async def getOjSourceCode(loginUrl):  
    width, height = 1400, 800  #网页宽高  
    browser = await pyp.launch(headless = False,  
                                userdataDir = "c:/tmp",  
                                args = [f'--window-size={width},{height}'])  
    page = await browser.newPage()  
    await antiAntiCrawler(page)  
    await page.setViewport({'width': width, 'height': height})  
    await page.goto(loginUrl)  
    #若手动登录, 则以下若干行可以去掉  
    element = await page.querySelector("#email")  #找到账户输入框  
    await element.type("XXXXXX@pku.edu.cn")  # 输入邮箱  
    element = await page.querySelector("#password")  #找到密码输入框  
    await element.type("XXXXXXXXXX")  # 输入密码
```



❏ pyppeteer 爬取 Openjudge 自己提交通过的所有程序源码

```
element = await page.querySelector("#main > form > div.user-login >
p:nth-child(2) > button")    #找到登录按钮
await element.click()        #点击登录按钮
#若手动登录, 则以上若干行可以去掉
await page.waitForSelector("#main>h2", timeout=30000)
#等待“正在进行的比赛...”标题出现
element = await page.querySelector("#userMenu>li:nth-child(2)>a")
#找“个人首页”链接
await element.click()        #点击个人首页链接
await page.waitForNavigation() #等新网页装入完毕
elements = await page.querySelectorAll(".result-right")
#找所有 Accepted 链接, 其有属性 class="result-right"
page2 = await browser.newPage() #新开一个页面 (标签)
await antiAntiCrawler(page2)
```



□ pyppeteer爬取Openjudge自己提交通过的所有程序源码

```
for element in elements[:2]: #只打印前两个程序
    obj = await element.getProperty("href") #获取href属性
    url = await obj.jsonValue()
    await page2.goto(url) #在新页面(标签)中装入新网页
    element = await page2.querySelector("pre") #查找pre tag
    obj = await element.getProperty("innerText") #取源代码
    text = await obj.jsonValue()
    print(text)
    print("-----")
await browser.close()

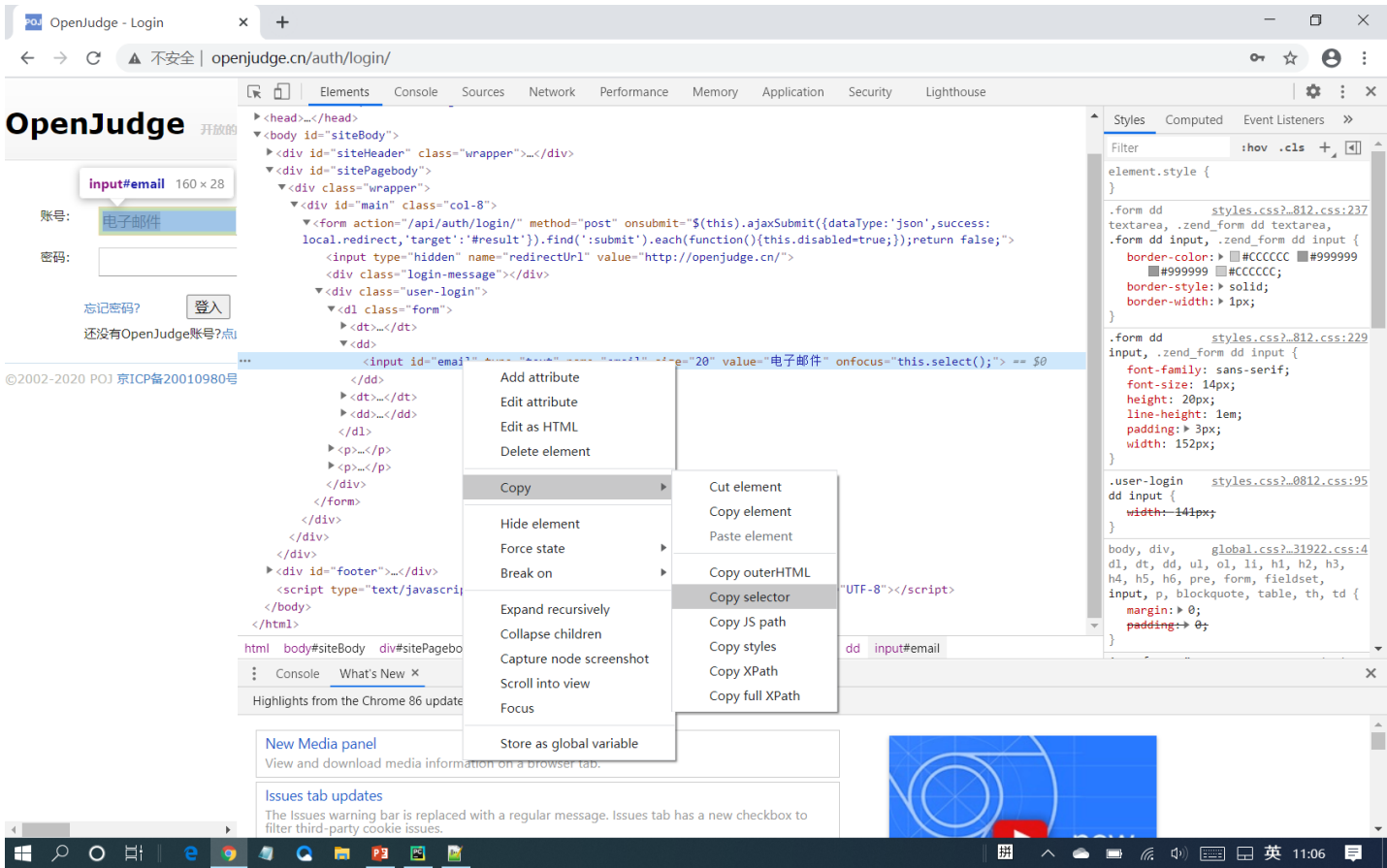
def main():
    url = "http://openjudge.cn/auth/login/"
    asyncio.get_event_loop().run_until_complete(getOjSourceCode(url))
main()
```



pyppeteer爬取Openjudge自己提交通过的所有程序源码

获取tag的
selector
得到selector:

#email



08 | 补充tips



□ 绝对网址和相对网址

绝对网址以 **http://** 或 **https://** 开头, 相对网址无这两种开头

如果当前网页网址是:

http://www.pku.edu.cn/education/index.htm

而该网页中有一个链接, 其中网址是相对的, 形如:

词典单词

则该链接的真实网址 (绝对网址) 是:

http://www.pku.edu.cn/education/dict/word.htm



□ 绝对网址和相对网址

使用requests库时, 获得当前网页网址:

```
r = requests.get("http://openjudge.cn")  
print(r.url) #>>http://openjudge.cn
```

或:

```
session = requests.session()  
r = session.get("http://openjudge.cn")  
print(r.url)
```

使用pypeteer库时, 获得当前网页网址:

```
browser = await pyp.launch(headless=False)  
page = await browser.newPage()  
await page.goto("http://openjudge.cn")  
print(page.url) #>>http://openjudge.cn
```



□ 反反爬技巧

连续的两个操作之间, 加入适当延时, 模拟人的动作, 避免因动作太快被识破

```
import time
```

```
time.sleep(2)    #暂停2秒, 啥也不做
```

也可以用`time.sleep(...)`来等待一段时间, 确保网页加载完成

