

第 3 关、BeautifulSoup 实践

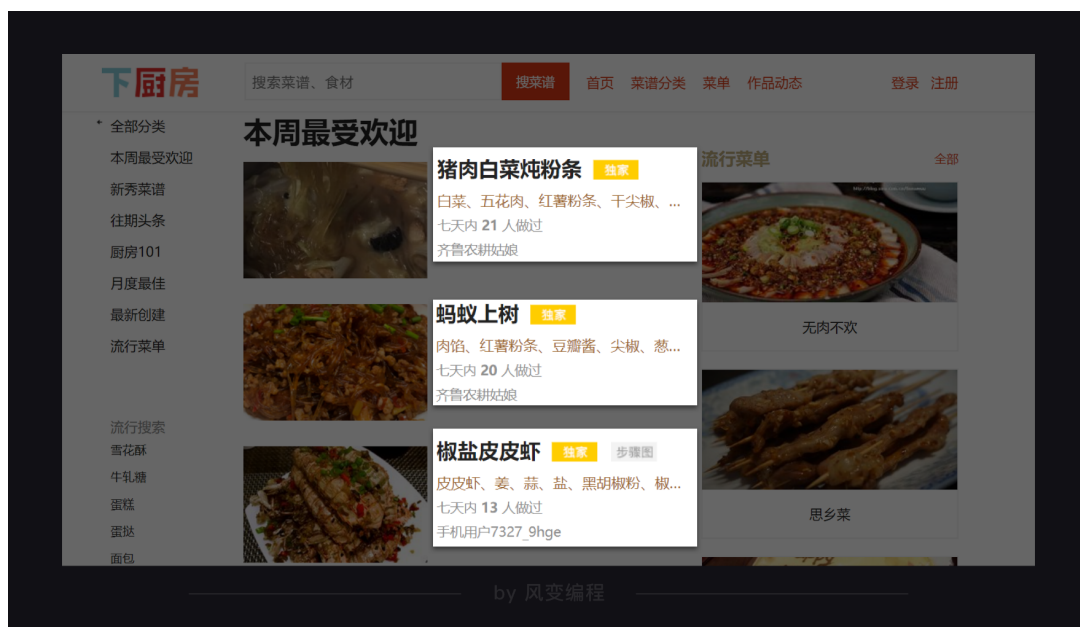
1、项目：解密吴氏私厨

1-1、确定目标

- (1) 目标网站：<http://www.xiachufang.com/explore/>；
- (2) 网站协议：<http://www.xiachufang.com/robots.txt>（目标网站 + robots.txt 可查看目标网站的页面爬取许可）；
- (3) 项目目标：爬取热门菜谱清单，内含：菜名、原材料、详细烹饪流程的URL。

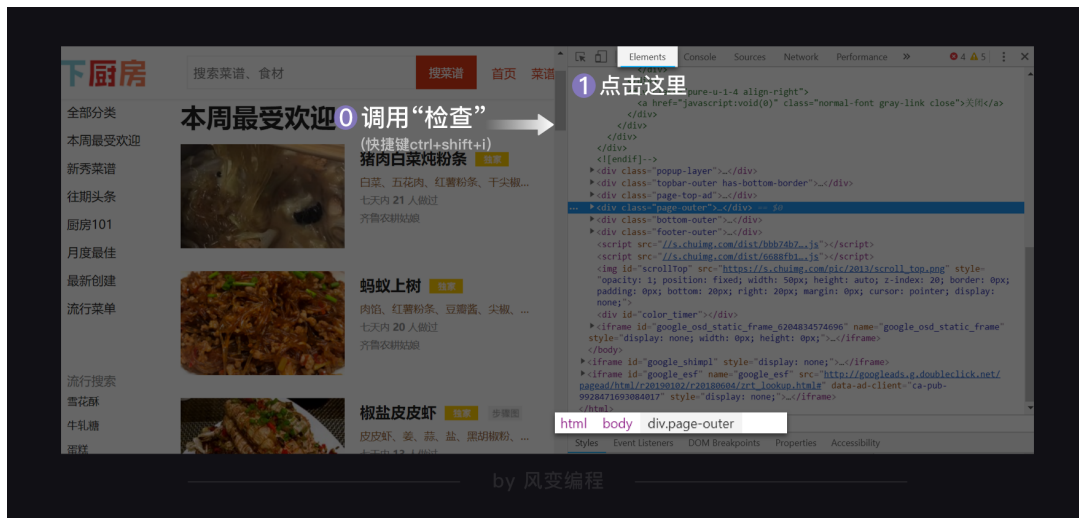
1-2、过程分析

- (1) 确定数据位置
 - 菜名、所需材料、和菜名所对应的详情页URL均在 html 页面上；
 - 获取数据用 `requests.get()`；
 - 解析数据用 `BeautifulSoup`。

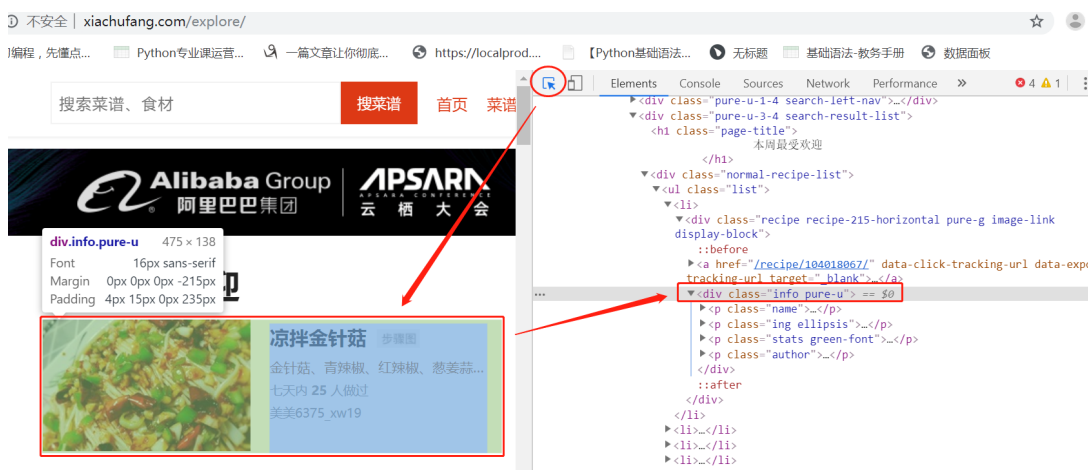


(2) 提取数据

- 【windows】：在网页的空白处点击右键，然后选择“检查”（快捷方式是ctrl+shift+i），再在 Elements 页面按 ctrl+f；【mac】：在网页的空白处点击右键，然后选择“检查”（快捷键 command + option + l(大写i)）；



- 点击【检查】页面左上角的“鼠标”按钮，再点击后右侧想要获取的内容可以定位到该内容对应的标签；



1-3、代码实现（一）

1-3-1、数据获取

requests.get() 获取数据，BeautifulSoup 解析数据。

```
1 import requests
2 # 引用requests库
3 from bs4 import BeautifulSoup
4 # 引用BeautifulSoup库
5
6 res_foods = requests.get('http://www.xiachufang.com/explore/')
7 # 获取数据
8 bs_foods = BeautifulSoup(res_foods.text, 'html.parser')
9 # 解析数据
10 print(bs_foods)
11 # 打印解析结果
```

1-3-2、提取最小父级标签

根据我们【过程分析】中所有菜谱的共同标签 class_='info pure-u'，我们用 find_all 获取所有菜谱（find_all 获取后返回的是一个列表），下面我们提取出第0个父级标签中的第0个 <a> 标签，并输出菜名和URL：

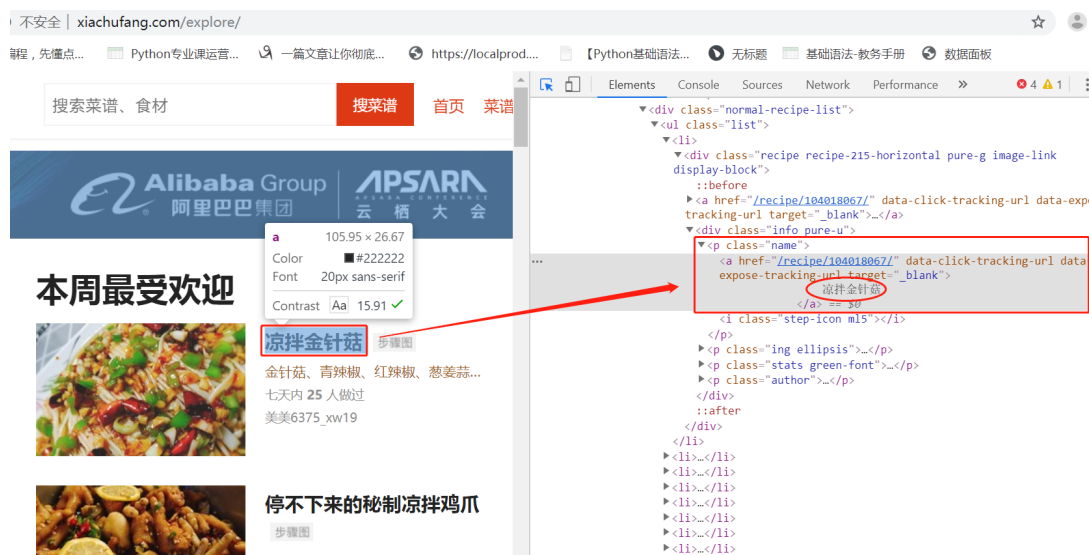
```

1 import requests
2 # 引用requests库
3 from bs4 import BeautifulSoup
4 # 引用BeautifulSoup库
5
6 res_foods = requests.get('http://www.xiachufang.com/explore/')
7 # 获取数据
8 bs_foods = BeautifulSoup(res_foods.text, 'html.parser')
9 # 解析数据
10 list_foods = bs_foods.find_all('div', class_='info pure-u')
11 # 查找最小父级标签
12 print(list_foods)
13 # 打印最小父级标签

```

1. 提取菜名

依旧是根据我们的内容定位我们的标签，可以找到菜名是在我们的标签 a 中，再用 text 取到该标签对应的菜名。



```

1 import requests
2 # 引用requests库
3 from bs4 import BeautifulSoup
4 # 引用BeautifulSoup库
5
6 res_foods = requests.get('http://www.xiachufang.com/explore/')
7 # 获取数据
8 bs_foods = BeautifulSoup(res_foods.text, 'html.parser')
9 # 解析数据
10 list_foods = bs_foods.find_all('div', class_='info pure-u')
11 # 查找最小父级标签
12
13 tag_a = list_foods[0].find('a')
14 # 提取第0个父级标签中的<a>标签
15 print(tag_a.text[17:-13])
16 # 输出菜名，使用[17:-13]切掉了多余的信息

```

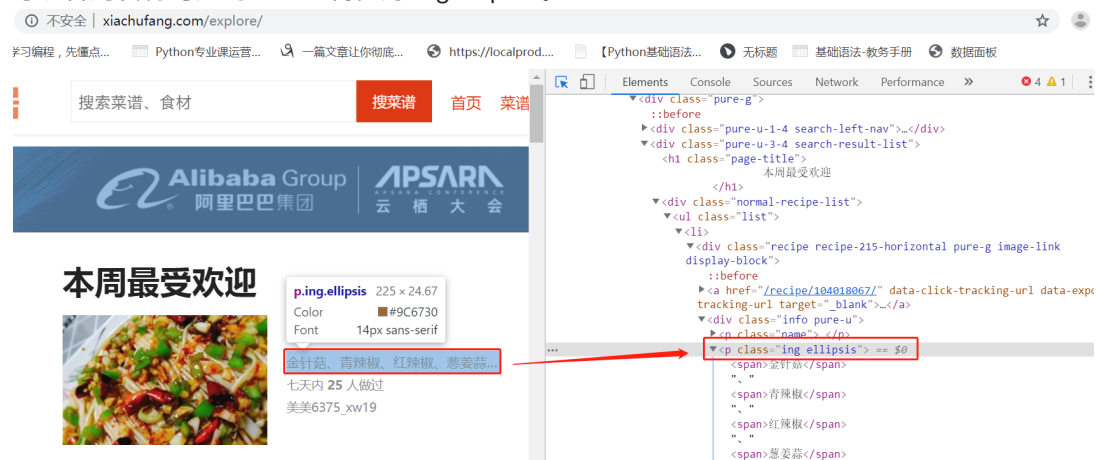
2. 提取 URL

我们发现在标签 a 后面的 href 有我们需要的链接，但是不完整，所以需要拼接后才能得到我们要的菜谱 URL。

```
1 import requests
2 # 引用requests库
3 from bs4 import BeautifulSoup
4 # 引用BeautifulSoup库
5
6 res_foods = requests.get('http://www.xiachufang.com/explore/')
7 # 获取数据
8 bs_foods = BeautifulSoup(res_foods.text, 'html.parser')
9 # 解析数据
10 list_foods = bs_foods.find_all('div', class_='info pure-u')
11 # 查找最小父级标签
12
13 tag_a = list_foods[0].find('a')
14 # 提取第0个父级标签中的<a>标签
15 print('http://www.xiachufang.com'+tag_a['href'])
16 # 拼接后输出URL
```

3. 提取食材

我们可以看到我们的食材是在 p 中，但是只靠这个是不够的，所以我们要精确取值，可以看到食材对应的 class 属性为 ing ellipsis。



```
1 import requests
2 # 引用requests库
3 from bs4 import BeautifulSoup
4 # 引用BeautifulSoup库
5
6 res_foods = requests.get('http://www.xiachufang.com/explore/')
7 # 获取数据
8 bs_foods = BeautifulSoup(res_foods.text, 'html.parser')
9 # 解析数据
10 list_foods = bs_foods.find_all('div', class_='info pure-u')
11 # 查找最小父级标签
12
13 tag_p = list_foods[0].find('p', class_='ing ellipsis')
14 # 提取第0个父级标签中的<p>标签
15 ingredients = tag_p.text[1:-1]
16 # 食材，使用[1:-1]切掉了多余的信息
```

```
17 print(ingredients)
18 # 打印食材
```

1-3-3、写循环，存列表

```
1 import requests
2 # 引用requests库
3 from bs4 import BeautifulSoup
4 # 引用BeautifulSoup库
5
6 res_foods = requests.get('http://www.xiachufang.com/explore/')
7 # 获取数据
8 bs_foods = BeautifulSoup(res_foods.text, 'html.parser')
9 # 解析数据
10 list_foods = bs_foods.find_all('div', class_='info pure-u')
11 # 查找最小父级标签
12
13 list_all = []
14 # 创建一个空列表，用于存储信息
15
16 for food in list_foods:
17
18     tag_a = food.find('a')
19     # 提取第0个父级标签中的<a>标签
20     name = tag_a.text[17:-13]
21     # 菜名，使用[17:-13]切掉了多余的信息
22     URL = 'http://www.xiachufang.com'+tag_a['href']
23     # 获取URL
24     tag_p = food.find('p', class_='ing ellipsis')
25     # 提取第0个父级标签中的<p>标签
26     ingredients = tag_p.text[1:-1]
27     # 食材，使用[1:-1]切掉了多余的信息
28     list_all.append([name, URL, ingredients])
29     # 将菜名、URL、食材，封装为列表，添加进list_all
30
31 print(list_all)
32 # 打印
```

1-4、代码实现（二）

创建一个空列表，启动循环，循环长度等于 <p> 标签的总数——我们可以借助 range(len()) 语法。

```
1 import requests
2 # 引用requests库
3 from bs4 import BeautifulSoup
4 # 引用BeautifulSoup库
5
6 res_foods = requests.get('http://www.xiachufang.com/explore/')
7 # 获取数据
8 bs_foods = BeautifulSoup(res_foods.text, 'html.parser')
```

```
9 # 解析数据
10
11 tag_name = bs_foods.find_all('p',class_='name')
12 # 查找包含菜名和URL的<p>标签
13 tag_ingredients = bs_foods.find_all('p',class_='ing ellipsis')
14 # 查找包含食材的<p>标签
15 list_all = []
16 # 创建一个空列表，用于存储信息
17 for x in range(len(tag_name)):
18 # 启动一个循环，次数等于菜名的数量
19     list_food = [tag_name[x].text[18:-14],tag_name[x].find('a')
20                  ['href'],tag_ingredients[x].text[1:-1]]
21     # 提取信息，封装为列表。注意此处[18:-14]切片和之前不同，是因为此处使用的是<p>
    标签，而之前是<a>
22     list_all.append(list_food)
23     # 将信息添加进list_all
24 print(list_all)
25 # 打印
```