

Lab1 报告

唐灵

519030910052

F1903002

2020 年 9 月 16 日

1 实验概览

总体来讲，本次实验可以分为四个子实验，即在利用 `python` 包的情况下，希望实现：

1. 通过 `url`，模拟浏览器访问服务器的过程，抓取网络页面源代码信息，当然，这次主要针对的是 `html` 文档。
2. 对页面源代码进行解析，生成 `dom` 树（将 `html` 文档表达为树结构），方便进行进一步的信息提取。
3. 通过专门处理 `dom` 树的包（`BeautifulSoup`），来限定并遍历访问希望访问的节点。
4. 通过正则表达式来匹配希望得到的特定信息，包括 `url`，文本信息等。
5. 最后一步不能算是实验，就是将得到的信息写入文件。

以上的四个小实验，所幸不用我们自己来实现，只需要调用基于 `python` 的不同功能的封装完整的包，很多时候就是一句话就能够解决的事情。

2 实验环境

本次实验的所有代码在电类工程导论 `c` 课程中在课程方统一给定的“`ee208`” `Docker` 容器中运行并实现，实验时间为 2020 年 09 月 11 日。

3 练习解决思路

3.1 练习一的解决思路

练习一的解决思路跟在之前提到的实验概览中提到的一样，有一点值得提出来一下，那就是实验限定了我们应当舍去图像的链接，我的做法是，在遍历匹配得到了页面中所有的指定标签下的 url 的情况下，在中间边进行遍历的过程中先判定末尾四位结尾是否为“.jpg”或者是“.png”，通过验证的 url 再加入到我们的 url 集合当中，以上为两种常见的图片格式，应当能够满足大部分的需求。

3.2 练习二的解决思路

练习二中依然和练习一是大致相同的思路，只不过这次需求相反，我们需要图片，这和我们希望匹配页面中的网页链接的步骤可以同步运行，只需要稍微改变正则表达式的值：为“`^http.*.(jpg|png)$`”即可（上标应该在最前边，而不是在字母‘h’上）其他的做法完全一致。

3.3 练习三的解决思路

针对问题三，与问题一二有两点不同：

1. 根据提示加入了消息头。
2. 其次是需要找到多组符合条件的信息，而非单一信息。

在按要求使用 `.add_headers()` 的方法，通过传递字典，添加了消息头信息之后。认真观察网页，审查网页元素，发现我们希望找到的元素都在同一类标签 `'a', {"class": "link - button"}` 之下，所以我们先通过 `bs` 找到这个标签，再进行遍历，每一次遍历过程处理这个节点下的子节点，提取信息，即可完成任务。

4 代码运行结果

4.1 练习一的运行结果

练习一结果储存在 `result1.txt` 之中，每一行代表一个 url，总共得到 34 条记录。

4.2 练习二的运行结果

练习二结果储存在 result1.txt 之中，每一行代表一个 url，最终发现图片链接仅有.png 格式，总过收获 10 条记录。

4.3 练习三的运行结果

练习三结果储存在 result1.txt 之中，每一行代表一组信息，总共得到了 30 组信息，部分截图如下 一段截图如下

```
6 https://pic1.zhimg.com/v2-9712000935e07e68015c61909d314.jpg?source=8673f162 网址 - 如何正确地喂 http://daily.zhihu.com/story/9727712
7 https://pic2.zhimg.com/v2-a8d171aaf4d53f096c75fd99c3c2.jpg?source=8673f162 第一次去空间站的话，怎么样操作会过去的样子? http://daily.zhihu.com/story/9727680
8 https://pic4.zhimg.com/v2-6a2a0c5ad2523a1164de8f08658483f.jpg?source=8673f162 空战时被击落的海军幸存的士兵会怎么样? http://daily.zhihu.com/story/9727683
9 https://pic4.zhimg.com/v2-31b75083320c2f43915c6020c431.jpg?source=8673f162 斯图亚特坦克和即时的影响? 没多大? http://daily.zhihu.com/story/9727693
10 https://pic3.zhimg.com/v2-f3ce8229f36c19995f8e081b04680.jpg?source=8673f162 为什么人类会长出这种奇怪的东西? http://daily.zhihu.com/story/9727711
11 https://pic4.zhimg.com/v2-ee0175c28c7ce100ea35b4485e8d399.jpg?source=8673f162 护肤品的活性成分添加越多，一定越好吗? http://daily.zhihu.com/story/9727703
12 https://pic3.zhimg.com/v2-676a7a67231cd9782ceadcced9f86d.jpg?source=8673f162 糟糕 - 如何正确地喂 http://daily.zhihu.com/story/9727679
13 https://pic4.zhimg.com/v2-472c586decd100bae5218500d8f.jpg?source=8673f162 有哪种比较奇怪的深海生物? http://daily.zhihu.com/story/9727644
14 https://pic2.zhimg.com/v2-3555b1fb6165320a721c4e539a3a10.jpg?source=8673f162 世界考古发现中有什么最让人吃惊的发现? http://daily.zhihu.com/story/9727658
15 https://pic1.zhimg.com/v2-9913076d1ac2c072115651eb682c2b00.jpg?source=8673f162 「银面人」是一种怎样的人体构造? http://daily.zhihu.com/story/9727678
16 https://pic1.zhimg.com/v2-82567e0430eb66fb204511861f0ba06.jpg?source=8673f162 为什么会饿? 不饿但不行啊? http://daily.zhihu.com/story/9727671
17 https://pic1.zhimg.com/v2-3531665fcd010e6e086231351af40b.jpg?source=8673f162 徒手劈砖需要多大力量? http://daily.zhihu.com/story/9727666
```

图 1: 练习三程序运行结果

5 分析与思考

本次实验进行了基本的网页爬取，没有什么很多好分析的，用包就完事儿了。

值得一提的是本来想通过打印加头和不加信息头的结果来看一下反爬机制，结果发现真的就返回了正确的结果。消息头的添加也让我感到比较迷惑，通过这个方法添加消息头，从形式上来讲不就只能传递 User-Agent 的信息吗？网上查了查，还是有其他的方法是通过传递字典来进行消息头的添加的。