

# Python Crawler

Python Crawler

Author: Yunqiu Xu

- 我以前学习爬虫的一些笔记, 主要基于Python核心编程和[静觅的爬虫教程](#)
- 

## Chapter 1 使用urllib/urllib2

### 1.1 构造一个最基本的爬取网页的代码

```
1. import urllib,urllib2
2. response = urllib2.urlopen("http://www.baidu.com")
3. print response.read()
```

- `urlopen(url, data, timeout):`
  - `data`: 访问url时要传送的数据,默认为 `None`, 1.3节则将用户名和密码作为该参数传输
  - `timeout`: 设定访问超时时间, 默认 `socket._GLOBAL_DEFAULT_TIMEOUT`

### 1.2 构造Request:

- 爬取过程: 构造request -- 服务器响应request -- 得到response

```
1. request=urllib2.Request("http://www.baidu.com")
2. response=urllib2.urlopen(request)
3. print response.read()
```

### 1.3 数据传输方式: POST&GET

- 动态网页需要动态传送参数给它(如需要输入用户名密码), 该网页做出对应的响应
  - GET 数据传输: 直接以链接形式访问, 连接中包含所有参数(不大安全, 但自己可以看到提交了什么内容)
-

```

1. #Get 数据传输
2. values={}
3. values['username']="username"
4. values['password']="password"
5. data=urllib.urlencode(values) #为dict编码，将编码后的dict作为data参数传输
6. url="url"
7. geturl=url+"?"+"data #url添加用户名密码等参数
8. request=urllib2.Request(geturl)
9. response=urllib2.urlopen(request)
10. print response.read()

```

- POST数据传输：不会在网址上显示所有的参数，查看提交的内容不大方便

```

1. #Post 数据传输
2. values={}
3. values['username']="username"
4. values['password']="password"
5. data=urllib.urlencode(values)
6. url="url"
7. request=urllib2.Request(url,data) #只有这里和Get不同
8. response=urllib2.urlopen(request)
9. print response.read()

```

## Chapter 2 `urllib/urllib2` 的更多用法

### 2.1 设置headers: 模拟浏览器工作

#### ▼ Request Headers [view source](#)

```

Accept: text/html,application/xhtml+xml,application/xml;q=0.9,image/webp,*/*;q=0.8
Accept-Encoding: gzip, deflate, sdch
Accept-Language: zh-CN,zh;q=0.8
Cache-Control: max-age=0
Connection: keep-alive
Cookie: laravel_session=eyJpdiI6Imw2c29PclB6cWlYK3NUTG1vbXNdNS3c9PSIsInZhbnHVLIjoisWNTZEtYTEZtVW84bHZNUUFcb3F5RCtuaDhHSklzXC9WSG82ZDRoWHdMTGZuWk5ZcTJobmhkR21WTXhHMjFkOXRicTBKVlh5OFJ4WjREblNFN2NFV0NRPT0iLCJtYWMiOiJhN2FhMTcyYzA3MzVjYjRkZWMyZm00TkYNDkZkbnVY0MDA2MDA2ODlLOGQyYTRkYWZhNzRkYTQ2MjgzZTUzIn0%3D; stat_fromWebUrl=; stat_ssaid=1474412240793; stat_uuid=1474290791895550452272; stat_isNew=1; _ga=GA1.2.1053141496.1474290342; gr_user_id=fcb6fe2f-9296-4ea1-b934-0fb098fb1e10; Hm_lvt_1320421bf6332884ce38374c23776351=1474290795; Hm_lpvt_1320421bf6332884ce38374c23776351=1474290795; QINGCLOUDELb=b2396366bd639e01fbd88e4a829ce868d562b81daa2a899c46d3a1e304c7eb2b|V9/ke|V9/KY
Host: wiki.jikexueyuan.com
Upgrade-Insecure-Requests: 1
User-Agent: Mozilla/5.0 (X11; Linux x86_64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/51.0.2704.84 Safari/537.36

```

- Request Headers 中包含许多信息，譬如agent代表请求身份，没有写入的话服务器可能会不响应
- 在原有的抓取程序中添加headers, 并将user\_agent加入这个dict用来让这个程序模拟浏览器，从而让服务器响应

```

1.     ...
2.     user_agent = "Mozilla/5.0 (X11; Linux x86_64) AppleWebKit/537.36
      (KHTML, like Gecko) Chrome/51.0.2704.84 Safari/537.36"
3.     headers={'User-Agent':user_agent}
4.     data=urllib.urlencode(data)
5.     request=urllib2.Request(url,data,headers=headers)
6.     response=urllib2.urlopen(request)
7.     page=response.read()

```

- 反盗链机制: 服务器识别header中的referer/host是不是其本身
  - 解决措施: 在header中添加referer

```

1.     headers['referer']="http://www.zhihu.com/"

```

## 2.2 设置代理

- 网站会检测某段时间某ip的访问次数，若访问次数过多会禁止访问
  - 解决措施: 每隔一段时间更换代理

```

1.     import urllib2
2.     enable_proxy = True
3.     proxy_handler = urllib2.ProxyHandler({"http" : 'http://some-
      proxy.com:8080'})
4.     null_proxy_handler = urllib2.ProxyHandler({})
5.     if enable_proxy:
6.         opener = urllib2.build_opener(proxy_handler)
7.     else:
8.         opener = urllib2.build_opener(null_proxy_handler)
9.     urllib2.install_opener(opener)

```

## 2.3 设置超时

- 停止访问一些反应过慢的页面
- 注意如果已经写明了data，直接传入超时时间即可,如10；反之，若data参数为空，需要特

别注明timeout=10

```
1. import urllib2
2. response = urllib2.urlopen('http://www.baidu.com',data, 10)
```

## Chapter 3 异常处理

- URLError

- 产生原因: 网络未连接; 连接不到特定服务器; 服务器不存在
- 捕获异常

```
1. import urllib2
2.
3. request = urllib2.Request('http://www.xxxx.com')
4. try:
5.     urllib2.urlopen(request)
6. except urllib2.URLError, e:
7.     print e.reason
```

- HTTPError:

- urlopen发送请求时, 服务器响应请求并返回response, 如果存在不能处理的内容则会产生HTTPError
- 注意HTTPError是URLError的一个子类, 子类异常优先被捕获, 捕获不到子类异常才会开始捕获父类

```
1. import urllib2
2.
3. req = urllib2.Request('http://blog.csdn.net/cqcre')
4. try:
5.     urllib2.urlopen(req)
6. except urllib2.HTTPError, e:
7.     print e.code
8. except urllib2.URLError, e:
9.     print e.reason
10. else:
11.     print "OK"
```

## Chapter 4 设置Cookie

- Cookie：某些网站为了辨别用户身份，进行session跟踪而储存在用户本地终端上的数据
- Opener：例如urlopen是一个特殊的opener,但我们使用Cookie时需要创建更一般的Opener
- Cookielib: 提供可存储cookie的对象以便与urllib2模块配合使用访问网络资源
  - 本模块内的CookieJar对象：捕获cookie并在后续连接请求时重新发送（如模拟登陆）
- 获取cookie并保存到变量:

```
1. import urllib2
2. import cookielib
3. #声明一个CookieJar对象实例来保存cookie
4. my_cookie = cookielib.CookieJar()
5.
6. #利用urllib2库的HTTPCookieProcessor对象来创建cookie处理器
7. handler=urllib2.HTTPCookieProcessor(my_cookie)
8.
9. #通过handler来构建opener
10. opener = urllib2.build_opener(handler)
11.
12. #此处的open方法同urllib2的urlopen方法，也可以传入request
13. response = opener.open('http://www.baidu.com')
14. for item in my_cookie:
15.     print 'Name = '+item.name
16.     print 'Value = '+item.value
17.
18. >>>
19. Name: BAIDUID
20. Value: 413C05FF289E8797034968BFADA0B1B8:FG=1
21. Name: BIDUPSID
22. Value: 413C05FF289E8797034968BFADA0B1B8
23. Name: H_PS_PSSID
24. Value: 1461_18240_21123_21193_21161_20928
25. Name: PSTM
26. Value: 1474296124
27. Name: BDSVRTM
28. Value: 0
29. Name: BD_HOME
30. Value: 0
```

- 获取cookie并保存到文件

```

1. import cookielib
2. import urllib2
3.
4. #设置保存cookie的文件，同级目录下的cookie.txt
5. filename = 'cookie.txt'
6. #声明一个MozillaCookieJar对象实例来保存cookie，之后写入文件
7. cookie = cookielib.MozillaCookieJar(filename)
8. #利用urllib2库的HTTPCookieProcessor对象来创建cookie处理器
9. handler = urllib2.HTTPCookieProcessor(cookie)
10. #通过handler来构建opener
11. opener = urllib2.build_opener(handler)
12. #创建一个请求，原理同urllib2的urlopen
13. response = opener.open("http://www.baidu.com")
14.
15. #保存cookie到文件
16. cookie.save(ignore_discard=True, ignore_expires=True)
17. # ignore_discard: 即使cookies将被丢弃也将它保存下来
18. # ignore_expire: 如果在该文件中cookies已经存在，则覆盖原文件写入

```

## ● 从文件中获取被保存的cookie进行访问

```

1. import cookielib
2. import urllib2
3.
4. #创建MozillaCookieJar实例对象
5. cookie = cookielib.MozillaCookieJar()
6. #从文件中读取cookie内容到变量
7. cookie.load('cookie.txt', ignore_discard=True, ignore_expires=True)
8. #创建请求的request
9. req = urllib2.Request("http://www.baidu.com")
10. #利用urllib2的build_opener方法创建一个opener
11. opener = urllib2.build_opener(urllib2.HTTPCookieProcessor(cookie))
12. response = opener.open(req)
13. print response.read()

```

## ● 使用cookie进行模拟登录

```

1. import urllib
2. import urllib2
3. import cookielib
4.
5. #声明一个MozillaCookieJar对象实例来保存cookie，之后写入文件
6. filename = 'cookie.txt'

```

```

7.  cookie = cookielib.MozillaCookieJar(filename)
8.  opener = urllib2.build_opener(urllib2.HTTPCookieProcessor(cookie))
9.  userdata = urllib.urlencode({
10.      'username':'Yunqiu Xu',
11.      'pwd':'12345678'
12.  })
13.
14.  #登录教务系统的URL
15.  urlname = 'urlname'
16.  #模拟登录, 并把cookie保存到变量
17.  result = opener.open(loginUrl,userdata)
18.  #保存cookie到cookie.txt中
19.  cookie.save(ignore_discard=True, ignore_expires=True)
20.
21.  #通过得到的cookie, 访问新网址
22.  urlname2 = 'urlname2'
23.  result = opener.open(urlname2)
24.  print result.read()

```

## Chapter 5 正则表达式

- 见*Python Notes Chapter 13*

## Chapter 6 `re` module的使用

- 见*Python Notes Chapter 13*
- 爬虫正则匹配基本框架

```

1.  pattern = re.compile("your regular expression here",re.S)
2.  items = re.findall(pattern,content)
3.  for item in items:
4.      print item[0],item[1],item[2],item[3],item[4]

```

## Chapter 7 一个简单的爬虫: 糗事百科段子抓取

- 抓取指定页数的段子(作者, 内容, 评分),排除其中有图片/视频的段子
- 源代码可见[我的github](#)

```

1.  #!/usr/bin/env python
2.  #coding: utf-8
3.
4.  import urllib
5.  import urllib2
6.  import re
7.
8.  #获取指定页码的内容
9.  def getPageFrom(url,current_page):
10.     referer = "http://www.qiushibaike.com/"
11.     user_agent = "Mozilla/5.0 (X11; Linux x86_64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/51.0.2704.84 Safari/537.36"
12.     headers = {'User-Agent': user_agent, 'referer': referer}
13.     new_url=url+'/page/'+str(current_page)
14.     request = urllib2.Request(new_url, headers=headers)
15.     response = urllib2.urlopen(request)
16.     content = response.read()
17.     return content
18.
19. #正则匹配：得到某一页的段子，返回段子的内容，作者和评分数，排除图片及视频
20. #得到作者
21. def get_author(content):
22.     pattern_author = re.compile('<h2>(.*?)</h2>',re.S)
23.     items_author = re.findall(pattern_author,content)
24.     for item in items_author:
25.         print item
26. #得到段子
27. def get_content(content):
28.     pattern_content = re.compile('<div class="content">[.\n]+<span>(.*?)</span>',re.S)
29.     items_content = re.findall(pattern_content,content)
30.     for item in items_content:
31.         print item
32. #得到评分数
33. def get_vote(content):
34.     pattern_vote = re.compile('<i class="number">(.*?)</i>',re.S)
35.     items_vote=re.findall(pattern_vote,content)
36.     for item in items_vote:
37.         print item
38. #综合以上结果，并排除有图片/视频的内容
39. def get_text_items(content):
40.     pattern_total = re.compile(
41.         '<h2>(.*?)</h2>.*?<div class="content">[.\n]+<span>(.*?)</span>(.*?)<i class="number">(.*?)</i>', re.S)
42.     items_total = re.findall(pattern_total, content)

```



```

43.     for item in items_total:
44.         if 'img' in item[2]:
45.             print "*****"
46.             print "带图片的段子，跳过!"
47.             continue
48.         else:
49.             print "*****"
50.             print "作者: " + item[0]
51.             print "内容: " + item[1]
52.             print "评分: " + item[3]
53.
54. #整合之前的函数
55. #输入起始页码和结束页码, 按页输出段子
56. def start(start_page,end_page):
57.     url = "http://www.qiushibaike.com/8hr"
58.     current_page=start_page
59.     while current_page<=end_page:
60.         print "现在开始打印第{}页的内容".format(current_page)
61.         content=getPageFrom(url,current_page)
62.         get_text_items(content)
63.         current_page+=1
64.         con=raw_input("继续吗?(Y/N) ")
65.         if con!='Y':
66.             print "停止打印"
67.             exit()
68.     print "完成全部输出!"
69.
70. if __name__=="__main__":
71.     print "糗事百科爬虫"
72.     start_page=int(raw_input("输入起始页码: "))
73.     end_page=int(raw_input("输入结束页码: "))
74.     start(start_page, end_page)

```

## Chapter 8 使用 BeautifulSoup

- BeautifulSoup提供一些简单的、python式的函数用来处理导航、搜索、修改分析树等功能
- BeautifulSoup自动将输入文档转换为Unicode编码，输出文档转换为utf-8编码, 不需要考虑编码方式，除非文档没有指定一个编码方式

## 8.1 创建BS对象

```
1. import urllib
2. import urllib2
3. from bs4 import BeautifulSoup
4.
5. #获取response的过程跟前面一样
6. request = urllib2.Request(new_url, headers=headers)
7. my_html=urllib2.urlopen(request)
8. #将response转化为BS对象
9. my_soup=BeautifulSoup(my_html)
10. print my_soup.prettify() #格式化输出
```

## 8.2 对象种类

### 8.2.1 Tag: HTML中的标签

```
1. print my_soup.title
2. >>> <title>糗事百科 - 超搞笑的原创糗事笑话分享社区</title>
3. print my_soup.a
4. >>> <a href="/"><h1>糗事百科</h1></a>
```

- Tag的两个重要属性: names, attrs
  - names为对象名, BS对象比较特殊输出[document],对于其他内部标签输出标签名
  - attrs为属性,得到字典
  - 属性可以进行赋值替换或删除

```
1. print my_soup.name
2. >>> [document]
3.
4. print my_soup.head.name
5. >>> head
6.
7. print my_soup.title.name
8. >>> title
9.
10. print my_soup.a.attrs
11. >>> {'href': '/'}
12.
13. print my_soup.div.attrs
14. >>> {'id': 'header', 'class': ['head']}
```

```
15.  
16.     del my_soup.div['class']
```

## 8.2.2 NavigableString: 可以遍历的字符串

- 得到标签内的内容: 省略了正则表达式

```
1.     print my_soup.title.string  
2.     >>> 糗事百科 - 超搞笑的原创糗事笑话分享社区  
3.  
4.     #若使用.strings则产生生成器对象
```

## 8.2.3 BeautifulSoup: 文档全部内容, 可以看做特殊的tag对象

```
1.     print my_soup.attrs  
2.     >>> {}  
3.  
4.     print my_soup.string  
5.     >>> None  
6.  
7.     print type(my_soup)  
8.     >>> <class 'bs4.BeautifulSoup'>
```

## 8.2.4 Comment: 特殊的NavigableString

- 建议在使用前先判断类型是否为Comment

```
1.     if type(soup.a.string)==bs4.element.Comment:  
2.         print soup.a.string
```

## 8.3 遍历文档树

- 子结点, 结点内容, 父结点, 兄弟结点, 前后结点

### 8.3.1 子结点:

- 直接子结点: `.contents`, `.children`
  - `.contents` 输出方式为列表
  - `.children` 输出方式为生成器
  - 使用for循环的话这两个的结果是一样的, 但 `.children` 不能直接输出列表

```
1.     for item in my_soup.head.contents:
```

```

2.     print item
3.
4.     for item in my_soup.head.children:
5.         print item

```

- 子孙结点: `.decendants`, 对所有tag的子孙结点进行递归循环

```

1.     for item in my_soup.descendants:
2.         print item

```

### 8.3.2 结点内容:

- `.string`
  - 如果一个标签里不存在标签,则返回当前标签中内容
  - 如果一个标签里只存在一个标签,则返回最里面标签的内容
  - 如果一个标签里存在多个标签,返回None

```

1.     print my_soup.head.string

```

- 多个内容: `.strings` / `.stripped_strings`
  - `.strings` 遍历获取多个内容
  - `.stripped_strings` 去除空行

```

1.     for item in my_soup.strings:
2.         print item
3.
4.     for item in my_soup.stripped.strings:
5.         print item

```

### 8.3.3 父结点:

- `.parent`

```

1.     print my_soup.div.parent.name
2.     >>> body

```

- 递归返回全部父结点: `.parents`

### 8.3.4 兄弟结点:

- `.next_sibling` / `.previous_sibling`

- 全部兄弟结点: `.next_siblings / .previous_siblings`

### 8.3.4 前后结点:

- 与兄弟结点不同: 前后结点针对所有结点, 不分层次
- `.next_element / .previous_element`
- `.next_elements / .previous_elements`

## 8.4 搜索文档树

### 8.4.1 `find_all( name , attrs , recursive , text , **kwargs )`

- 遍历当前tag的所有子结点
- `name`: 查找所有名字为name的tag

```

1.  #传入字符串
2.  my_soup.find_all('a')
3.  my_soup.find_all('div')
4.
5.  #传入正则表达式
6.  pattern=re.compile('regular expression')
7.  my_soup.find_all(pattern)
8.
9.  #传入列表: 只要匹配列表中某一个函数就会返回
10. my_soup.find_all(['a','div'])
11.
12. #传入True: 返回所有标签, 但不会返回字符串结点
13. my_soup.find_all(True)
14. my_soup.find_all() #等价
15.
16. #传入方法: 定义一个方法接受当前元素, 返回True则表示匹配成功
17. #有些类似filter函数
18. def has_div(item):
19.     return 'div' in item
20. my_soup.find_all(has_div)

```

- 关键字参数

```

1.  my_soup.find_all("span",class_="sister")
2.  my_soup.find_all("span",id="Yunqiu Xu")

```

- **attrs参数:** 一些包含特殊属性的tag无法搜索, 可通过attrs参数转换为字典参数

```
1. #HTML5中的data-*属性不能直接搜索
2. data_soup = BeautifulSoup('<div data-foo="value">foo!</div>')
3. data_soup.find_all(data-foo="value")
4. >>> SyntaxError: keyword can't be an expression
5.
6. #使用attrs转换
7. data_soup.find_all(attrs={"data-foo": "value"})
8. >>> [<div data-foo="value">foo!</div>]
```

- **text参数:** 搜索文档中的字符串, 参数可选值同name

```
1. my_soup.find_all(text=True)
2. my_soup.find_all(text='sb')
```

- **limit参数:** 限制返回结果数量

```
1. my_soup.findall(limit=10)
```

- **recursive参数:**
  - 默认为True: 遍历全部子孙结点
  - False: 仅仅返回直接子结点

```
1. len(my_soup.find_all('div'))
2. >>> 215
3. len(my_soup.find_all('div', recursive=False))
4. >>> 0
```

## 8.4.2 其他搜索函数

```
1. find() #参数同find_all(), 但是仅仅返回第一个匹配的结果
2.
3. #返回当前匹配结点的父辈结点
4. find_parent()
5. find_parents()
6.
7. #返回当前匹配结点的兄弟结点
8. find_next_sibling()
9. find_next_siblings()
10. find_previous_sibling()
11. find_previous_siblings()
```

```
12.
13.     #返回当前匹配结点的前后结点
14.     find_next()
15.     find_all_next()
16.     find_previous()
17.     find_all_previous()
```

## 8.5 CSS选择器 `soup.select()`

- 标签名查找

```
1.     my_soup.select('title')
```

- 同理, 还有类名查找/id名查找, 或者多种查找方式的组合查找(使用空格分割)

```
1.     my_soup.select('head title')
```

- 属性查找: 属性与标签属于同一结点, 中间不可加空格

```
1.     soup.select('a[href="http://example.com/elsie"]')
```

## Chapter 9 `requests` module

- urllib2提供了爬虫所需的大部分功能, 但其过于复杂, requests库用于简化这些功能

### 9.1 发出请求

```
1.     import requests
2.
3.     url='http://www.qiushibaike.com/'
4.     r1=requests.get(url)
5.     r2=requests.post(url,data={"username":"password"})
6.     r3=requests.put(url,data={"username":"password"})
7.     r4=requests.head(url)
8.
9.     #传递参数
10.    data={"username":"password"}
```

```
11. r5=requests.get(url,params=data)
12.
13. #设置超时
14. r6=requests.get(url,timeout=0.001)
```

## 9.2 响应内容

```
1. r.text #相当于之前的response.read()
2. r.encoding #获取网页编码
3. >>> 'UTF-8'
4. r.content #二进制响应
5. r.json() #json响应
6.
7. #获取响应状态
8. r.status_code
9. >>> 405
10.
11. #响应headers
12. r.headers
```

## 9.3 模拟登陆

```
1. # 设置headers
2. referer = "http://www.qiushibaike.com/"
3. user_agent = "Mozilla/5.0 (X11; Linux x86_64) AppleWebKit/537.36
4. (KHTML, like Gecko) Chrome/51.0.2704.84 Safari/537.36"
5. headers = {'User-Agent': user_agent, 'referer': referer}
6. r=requests.get(url,headers=headers)
7.
8. # 传递一个dict到data参数处
9. my_data={'key1':'val1','key2':'val2'}
10. r=requests.post(url,data=my_data)
11. # 也可以传递到json参数处
12. r=requests.post(url,json=my_data)
13. # 从文件中提取数据传递到file参数
14. my_file={'file':open(path,'rb')}
15. r=requests.post(url,file=my_file)
```

## 9.4 获取并传递cookies



- 获取cookies

```
1. r=requests.get(url,headers=headers)
2. my_cookies=r.cookies
```

- 将自己的cookies发送到服务器

```
1. my_cookies=={'cookies':'xxx'}
2. r=requests.get(url,cookies=my_cookies)
```

## 9.5 查询历史及重定向

```
1. r.history #查询历史
2. r=requests.get(url,allow_redirects=False) #禁止重定向
```

## 9.6 高级用法

- 可参考[官方文档](#)
- To be continued