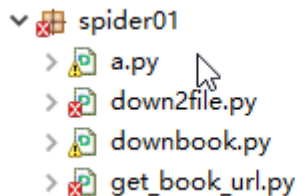


程序目的：根据笔趣阁网站的排行榜：<https://www.biqukan.com/paihangbang/>，获取每部小说的目录页，然后逐页爬取小说，并写入文件。

程序未实现多线程，目前爬取速度缓慢。

目录：



a.py 为主程序，

down2file.py: 调用downbook.py的类，逐页下载写入文件

downbook.py: 据提供的网站url（用于程序中url的补全）和小说的目录页url，实现了内容的获取（函数）和文件写入（函数）

get_book_url.py:根据排行榜页获取每个小说的目录页

a.py

```
1  # -*- coding=utf-8 -*-
2  '''
3  程序说明：爬虫小说网站的主程序
4  '''
5  from spider01 import down2file,get_book_url
6  from multiprocessing import Pool
7
8
9  myurl='https://www.biqukan.com/paihangbang/' #这个是排行榜那个网址
10 myup_url='http://www.biqukan.com'#这个是网站url，用于后面程序的生成url时补全头部
11 get01=get_book_url #实例化
12 mybookurllist,mybooknamelist=get01.getlist(url=myurl,up_url=myup_url) #获取小说目录页url的列表、小说名称
13 book=down2file #实例化
14
15 for i in range(len(mybookurllist)):
16     book.down_file(myup_url,mybookurllist[i],mybooknamelist[i]) #逐本下载
17     print("所有下载都已完成")
18
```

down2file.py

```

1 from spider01 import downbook
2 import sys ,random
3 from time import sleep
4 '''
5 本程序内类主要实现调用downbook的downloader类，
6 逐页下载写入相应文件
7 '''
8
9
10 class down_file:
11     def __init__(self,server,bookurl,bookname):
12         self.server=server
13         self.bookurl=bookurl
14         self.bookname=bookname
15         dl = downbook.downloader(server=self.server,bookurl=self.bookurl)
16         dl.get_download_url()
17         print(' 《 %s》 开始下载: '%self.bookname)
18         #for i in range(2): #减少每个小说的下载章数，用于测试程序的运行
19         for i in range(dl.nums):
20             sleep(10+random.random())#添加random时间
21             dl.writer(dl.names[i], '%s.txt'%self.bookname,
22 dl.get_contents(dl.urls[i]))
23             sys.stdout.write(" %s已下载:%.3f%%" % (self.bookname,float(i/dl.nums)) +
24 '\n')
25             sys.stdout.flush()
26         print(' 《%s》 下载完成'%self.bookname)

```

downbook.py

```

1 # -*- coding:UTF-8 -*-
2 from bs4 import BeautifulSoup
3 import requests
4
5 """
6 类说明:根据提供的网站url（用于程序中url的补全）和小说的目录页url，实现了内容的获取（函数）和文件写入（函数）
7 Parameters:
8     server: 网站url（用于程序中url的补全）
9     bookurl: 小说的目录页url

```

```

10 Returns:
11 无
12 Modify:
13 2018-12-4
14 """
15 class downloader(object):
16
17     def __init__(self, server, bookurl):
18         self.server = server
19         self.target = bookurl
20         self.names = [] #存放章节名
21         self.urls = [] #存放章节链接
22         self.nums = 0 #章节数
23
24     """
25     函数说明: 获取下载链接
26     Parameters:
27     无
28     Returns:
29     无
30     Modify:
31     2018-12-4
32     """
33     def get_download_url(self):
34         req = requests.get(url = self.target)
35         html = req.text
36         div_bf = BeautifulSoup(html, features='lxml')
37         div = div_bf.find_all('div', class_ = 'listmain')
38         a_bf = BeautifulSoup(str(div[0]), features='lxml')
39         a = a_bf.find_all('a')
40         self.nums = len(a[12:]) #剔除不必要的章节, 并统计章节数
41         for each in a[12:]:
42             self.names.append(each.string)
43             self.urls.append(self.server + each.get('href'))
44
45     """
46     函数说明: 获取章节内容
47     Parameters:
48     target - 下载连接(string)
49     Returns:

```

```

50  texts - 章节内容(string)
51  Modify:
52  2018-12-4
53  """
54  def get_contents(self, target):
55      req = requests.get(url = target)
56      html = req.text
57      bf = BeautifulSoup(html, features='lxml')
58      texts = bf.find_all('div', class_ = 'showtxt')
59      texts = texts[0].text.replace('\xa0'*8, '\n\n')
60      return texts
61
62  """
63  函数说明:将爬取的文章内容写入文件
64  Parameters:
65  name - 章节名称(string)
66  path - 当前路径下,小说保存名称(string)
67  text - 章节内容(string)
68  Returns:
69  无
70  Modify:
71  2018-12-4
72  """
73  def writer(self, name, path, text):
74      write_flag = True
75      with open(path, 'a', encoding='utf-8') as f:
76          f.write(name + '\n')
77          f.writelines(text)
78          f.write('\n\n')
79
80

```

get_book_url.py

```

1  #-*- coding:UTF-8 -*-
2  from bs4 import BeautifulSoup
3  import requests
4
5
6  '''

```

```

7 类说明:获取网页上单个小说的url（目录页）列表
8  '''
9  '''
10 class Getbookurl:
11     def __init__(self):
12         pass
13     '''
14
15     '''
16     函数说明: 获取小说的目录页
17     Parameters:
18     url: 小说网页
19     up_url: 上级目录的url
20     Returns:无
21     Modify:2018-12-4
22     '''
23     def getlist(url,up_url):
24         print("开始整理小说目录页的url列表,请稍等!")
25         req=requests.get(url)
26         html=req.text
27         div_bf=BeautifulSoup(html,features='lxml')
28         div=div_bf.find_all('div', class_='block bd')
29         bookurls=[] #小说目录页列表
30         booknames=[]#小说名称列表
31         for i in range(len(div)):
32             a_bf=BeautifulSoup(str(div[i]),features='lxml')
33             a=a_bf.find_all('a')
34             booknums=len(a)
35             for each in a:
36                 bookurls.append(up_url+each.get('href'))
37                 booknames.append(each.string)
38             print('小说目录页的url列表整理完成!',bookurls)
39             print('所有需要下载的小说',booknames)
40         return bookurls,booknames
41

```

运行a.py,最终实现每个小说的下载。

目前程序还是单线程运行，为了规避网站“反爬虫”，设置了随机sleep时间。这样效率下降了很多。