

DAY02

Day01回顾

请求模块(*urllib.request*)

```
1 req = request.Request(url,headers=headers)
2 res = request.urlopen(req)
3 html = res.read().decode('utf-8')
```

编码模块(*urllib.parse*)

```
1 1、urlencode({dict})
2     urlencode({'wd':'美女','pn':'20'})
3     编码后 : 'wd=%E8%D5XXX&pn=20'
4
5 2、quote(string)
6     quote('织女')
7     编码后 : '%D3%F5XXX'
8
9 3、unquote('%D3%F5XXX')
```

解析模块(*re*)

■ 使用流程

```
1 pattern = re.compile('正则表达式',re.S)
2 r_list = pattern.findall(html)
```

■ 贪婪匹配和非贪婪匹配

```
1 贪婪匹配(默认) : .*
2 非贪婪匹配      : .*?
```

■ 正则表达式分组

- 1 1、想要什么内容在正则表达式中加()
- 2 2、多个分组,先按整体正则匹配,然后再提取()中数据。结果: [(),(),(),(),()]

抓取步骤

- 1 1、确定所抓取数据在响应中是否存在 (右键 - 查看网页源码 - 搜索关键字)
- 2 2、数据存在: 查看URL地址规律
- 3 3、写正则表达式,来匹配数据
- 4 4、程序结构
 - 5 1、使用随机User-Agent
 - 6 2、每爬取1个页面后随机休眠一段时间

```
1 # 程序结构
2 class xxxSpider(object):
3     def __init__(self):
4         # 定义常用变量,url,headers及计数等
5
6     def get_html(self):
7         # 获取响应内容函数,使用随机User-Agent
8
9     def parse_html(self):
10        # 使用正则表达式来解析页面,提取数据
11
12    def write_html(self):
13        # 将提取的数据按要求保存, csv、MySQL数据库等
14
15    def main(self):
16        # 主函数,用来控制整体逻辑
17
18 if __name__ == '__main__':
19     # 程序开始运行时间戳
20     start = time.time()
21     spider = xxxSpider()
22     spider.main()
23     # 程序运行结束时间戳
24     end = time.time()
25     print('执行时间:%.2f' % (end-start))
```

Day02笔记

作业讲解

作业1 - 正则分组练习

页面结构如下:

```

1  <div class="animal">
2      <p class="name">
3          <a title="Tiger"></a>
4      </p>
5      <p class="content">
6          Two tigers two tigers run fast
7      </p>
8  </div>
9
10 <div class="animal">
11     <p class="name">
12         <a title="Rabbit"></a>
13     </p>
14
15     <p class="content">
16         Small white rabbit white and white
17     </p>
18 </div>

```

从以上html代码结构中完成如下内容信息的提取：

```

1  # 问题1
2  [('Tiger', 'Two...'), ('Rabbit', 'Small..')]
3  # 问题2
4  动物名称 : Tiger
5  动物描述 : Two tigers two tigers run fast
6  *****
7  动物名称 : Rabbit
8  动物描述 : Small white rabbit white and white

```

代码实现

```

1  import re
2
3  html = '''
4  <div class="animal">
5      <p class="name">
6          <a title="Tiger"></a>
7      </p>
8      <p class="content">
9          Two tigers two tigers run fast
10     </p>
11 </div>
12
13 <div class="animal">
14     <p class="name">
15         <a title="Rabbit"></a>
16     </p>
17
18     <p class="content">
19         Small white rabbit white and white
20     </p>
21 </div>
22 '''

```

```

23 pattern = re.compile(r'<div class="animal">.*?<a title="(.*?)".*?content">(.*?)</p>',re.S)
24 r_list = pattern.findall(html)
25 # 问题1
26 if r_list:
27     print(r_list)
28 # r_list: [('Tiger','\n\t\t Two tigers xxx'),()]
29 # 问题2
30 if r_list:
31     for r in r_list:
32         print('动物名称:',r[0].strip())
33         print('动物描述:',r[1].strip())
34         print('*' * 50)
35 else:
36     print('未匹配到数据')

```

猫眼电影top100抓取案例

- 1 猫眼电影 - 榜单 - top100榜
- 2 电影名称、主演、上映时间

数据抓取实现

■ 1、确定响应内容中是否存在所需数据

- 1 右键 - 查看网页源代码 - 搜索关键字 - 存在！！

■ 2、找URL规律

- 1 第1页: <https://maoyan.com/board/4?offset=0>
- 2 第2页: <https://maoyan.com/board/4?offset=10>
- 3 第n页: `offset=(n-1)*10`

■ 3、正则表达式

- 1 `<div class="movie-item-info">.*?title="(.*?)".*?class="star">(.*?)</p>.*?releasetime">(.*?)</p>`

■ 4、编写程序框架，完善程序

```

1 from urllib import request
2 import re
3 import time
4 import random
5 from useragents import ua_list
6
7 class MaoyanSpider(object):
8     def __init__(self):
9         self.url = 'https://maoyan.com/board/4?offset={}'
10        # 计数
11        self.num = 0

```

```

12
13 # 获取
14 def get_html(self,url):
15     headers = {
16         'User-Agent' : random.choice(ua_list)
17     }
18     req = request.Request(url=url,headers=headers)
19     res = request.urlopen(req)
20     html = res.read().decode('utf-8')
21     # 直接调用解析函数
22     self.parse_html(html)
23
24 # 解析
25 def parse_html(self,html):
26     re_bds = r'<div class="movie-item-info">.*?title="(.*?)".*?class="star">(.*?)</p>.*?
releasetime">(.*?)</p>'
27     pattern = re.compile(re_bds,re.S)
28     # film_list: [('霸王别姬','张国荣','1993'),()]
29     film_list = pattern.findall(html)
30     # 直接调用写入函数
31     self.write_html(film_list)
32
33 def write_html(self,film_list):
34     item = {}
35     for film in film_list:
36         item['name'] = film[0].strip()
37         item['star'] = film[1].strip()
38         item['time'] = film[2].strip()[5:15]
39         print(item)
40
41         self.num += 1
42
43 def main(self):
44     for offset in range(0,31,10):
45         url = self.url.format(offset)
46         self.get_html(url)
47         time.sleep(random.randint(1,2))
48         print('共抓取数据:',self.num)
49
50 if __name__ == '__main__':
51     start = time.time()
52     spider = MaoyanSpider()
53     spider.main()
54     end = time.time()
55     print('执行时间:%.2f' % (end-start))

```

数据持久化存储

数据持久化存储 - csv文件

▪ 作用

```
1 将爬取的数据存放到本地的csv文件中
```

▪ 使用流程

```
1 1、导入模块
2 2、打开csv文件
3 3、初始化写入对象
4 4、写入数据(参数为列表)
5 import csv
6
7 with open('film.csv','w') as f:
8     writer = csv.writer(f)
9     writer.writerow([])
```

▪ 示例代码

创建 test.csv 文件，在文件中写入数据

```
1 # 单行写入 (writerow([]))
2 import csv
3 with open('test.csv','w',newline='') as f:
4     writer = csv.writer(f)
5     writer.writerow(['步惊云','36'])
6     writer.writerow(['超哥哥','25'])
7
8 # 多行写入(writerows([(),(),())])
9 import csv
10 with open('test.csv','w',newline='') as f:
11     writer = csv.writer(f)
12     writer.writerows([('聂风','36'),('秦霜','25'),('孔慈','30')])
```

▪ 练习

猫眼电影数据存入本地 maoyanfilm.csv 文件 - 使用writerow()方法实现

```
1 # 存入csv文件 - writerow()
2 def write_html(self,film_list):
3     with open('film.csv','a') as f:
4         # 初始化写入对象,注意参数f别忘了
5         writer = csv.writer(f)
6         for film in film_list:
7             L = [
8                 film[0].strip(),
9                 film[1].strip(),
10                film[2].strip()[5:15]
11            ]
12            # writerow()参数为列表
13            writer.writerow(L)
```

思考：使用 writerows()方法实现？

```
1 # 存入csv文件 - writerows()
```

```

2 def write_html(self, film_list):
3     L = []
4     with open('film.csv', 'a') as f:
5         # 初始化写入对象, 注意参数f别忘了
6         writer = csv.writer(f)
7         for film in film_list:
8             t = (
9                 film[0].strip(),
10                film[1].strip(),
11                film[2].strip()[5:15]
12            )
13            L.append(t)
14            # writerows()参数为列表
15            writer.writerows(L)

```

数据持久化存储 - MySQL 数据库

1、在数据库中建库建表

```

1 # 连接到mysql数据库
2 mysql -h127.0.0.1 -uroot -p123456
3 # 建库建表
4 create database maoyandb charset utf8;
5 use maoyandb;
6 create table filmtab(
7     name varchar(100),
8     star varchar(300),
9     time varchar(50)
10 )charset=utf8;

```

2、回顾pymysql基本使用

```

1 import pymysql
2
3 # 创建2个对象
4 db = pymysql.connect('localhost', 'root', '123456', 'maoyandb', charset='utf8')
5 cursor = db.cursor()
6
7 # 执行SQL命令并提交到数据库执行
8 # execute()方法第二个参数为列表传参补位
9 ins = 'insert into filmtab values(%s,%s,%s)'
10 cursor.execute(ins, ['霸王别姬', '张国荣', '1993'])
11 db.commit()
12
13 # 关闭
14 cursor.close()
15 db.close()

```

来试试高效的executemany()方法?

```

1 import pymysql

```

```

2
3 # 创建2个对象
4 db = pymysql.connect('192.168.153.137','tiger','123456','maoyandb',charset='utf8')
5 cursor = db.cursor()
6
7 # 抓取的数据
8 film_list = [('月光宝盒','周星驰','1994'),('大圣娶亲','周星驰','1994')]
9
10 # 执行SQL命令并提交到数据库执行
11 # execute()方法第二个参数为列表传参补位
12 cursor.executemany('insert into filmtab values(%s,%s,%s)',film_list)
13 db.commit()
14
15 # 关闭
16 cursor.close()
17 db.close()

```

■ 3、将电影信息存入MySQL数据库（尽量使用executemany方法）

```

1 # mysql - executemany([()],(),())
2 def write_html(self, film_list):
3     L = []
4     ins = 'insert into filmtab values(%s,%s,%s)'
5     for film in film_list:
6         t = (
7             film[0].strip(),
8             film[1].strip(),
9             film[2].strip()[5:15]
10        )
11        L.append(t)
12
13        self.cursor.executemany(ins, L)
14        # 千万别忘了提交到数据库执行
15        self.db.commit()

```

■ 4、做个SQL查询

```

1 1、查询20年以前的电影的名字和上映时间
2     select name,time from filmtab where time<(now()-interval 20 year);
3 2、查询1990-2000年的电影名字和上映时间
4     select name,time from filmtab where time>='1990-01-01' and time<='2000-12-31';

```

数据持久化存储 - MongoDB数据库

pymongo操作mongodb数据库


```

1 import pymongo
2
3 # 1.数据库连接对象
4 conn=pymongo.MongoClient('localhost',27017)
5 # 2.库对象
6 db = conn['库名']
7 # 3.集合对象
8 myset = db['集合名']
9 # 4.插入数据
10 myset.insert_one({字典})

```

思考

- 1、能否到电影详情页把评论抓取下来?
- 2、能否到电影详情页把电影图片抓取下来? - 并按照电影名称分别创建文件夹

代码实现

```

1 from urllib import request
2 import re
3 import time
4 import random
5 from useragents import ua_list
6 import os
7
8 class MaoyanSpider(object):
9     def __init__(self):
10         self.url = 'https://maoyan.com/board/4?offset={}'
11         # 计数
12         self.num = 0
13
14     # 获取
15     def get_html(self,url):
16         headers = {
17             'User-Agent' : random.choice(ua_list)
18         }
19         req = request.Request(url=url,headers=headers)
20         res = request.urlopen(req)
21         html = res.read().decode('utf-8','ignore')
22
23         return html
24
25     def re_func(self,re_bds,html):
26         pattern = re.compile(re_bds,re.S)
27         r_list = pattern.findall(html)
28
29         return r_list
30
31     # 解析
32     def parse_html(self,url):
33         re_bds = r'<div class="movie-item-info">.*?href="(.*?)".*?title="(.*?)".*?class="star">
34         (.*?)</p>.*?releasetime">(.*?)</p>'
35         # html获取 + re解析
36         html = self.get_html(url)

```

```

36     film_list = self.re_func(re_bds,html)
37     # 直接调用写入函数
38     self.write_html(film_list)
39
40 def write_html(self,film_list):
41     film_dict = {}
42     for film in film_list:
43         film_dict['name'] = film[1].strip()
44         film_dict['star'] = film[2].strip()
45         film_dict['time'] = film[3].strip()[5:15]
46         two_url = 'https://maoyan.com{}'.format(film[0].strip())
47         film_dict['comment'] = self.get_comment(two_url)
48         self.save_image(two_url,film)
49         print(film_dict)
50
51         self.num += 1
52
53 def get_comment(self,two_url):
54     # 获取 + 解析
55     html = self.get_html(two_url)
56     re_bds = r'<div class="comment-content">(.*?)</div>'
57     comment_list = self.re_func(re_bds,html)
58
59     return comment_list
60
61 def save_image(self,two_url,film):
62     re_bds = r'<div class="img.*?"></div>'
63     html = self.get_html(two_url)
64     img_link_list = self.re_func(re_bds,html)
65
66
67     for img_link in img_link_list:
68         req = request.Request(img_link)
69         res = request.urlopen(req)
70         html = res.read()
71         # 处理文件名
72         directory = 'D:\\猫眼\\{\\}\\'.format(film[1].strip())
73         if not os.path.exists(directory):
74             os.makedirs(directory)
75
76         filename=directory + img_link.split('/')[ -1].split('@')[0]
77         with open(filename,'wb') as f:
78             f.write(html)
79
80
81 # 入口函数
82 def run(self):
83     for offset in range(0,31,10):
84         url = self.url.format(offset)
85         self.parse_html(url)
86         time.sleep(random.randint(1,2))
87         print('共抓取数据:',self.num)
88
89 if __name__ == '__main__':
90     start = time.time()
91     spider = MaoyanSpider()
92     spider.run()

```

```
93     end = time.time()
94     print('执行时间:%.2f' % (end-start))
```

电影天堂二级页面抓取案例

领取任务

```
1  # 地址
2  电影天堂 - 2019年新片精品 - 更多
3  # 目标
4  电影名称、下载链接
5
6  # 分析
7  *****一级页面需抓取*****
8      1、电影详情页链接
9
10 *****二级页面需抓取*****
11     1、电影名称
12     2、电影下载链接
```

实现步骤

- 1、确定响应内容中是否存在所需抓取数据
- 2、找URL规律

```
1 第1页 : https://www.dytt8.net/html/gndy/dyzz/list_23_1.html
2 第2页 : https://www.dytt8.net/html/gndy/dyzz/list_23_2.html
3 第n页 : https://www.dytt8.net/html/gndy/dyzz/list_23_n.html
```

- 3、写正则表达式

```
1 1、一级页面正则表达式
2     <table width="100%".*?<td width="5%".*?<a href="(.*?)".*?<u link">.*?</table>
3 2、二级页面正则表达式
4     <div class="title_all"><h1><font color=#07519a>(.*?)</font></h1></div>.*?<td style="WORD-
    WRAP.*?>.*?>(.*?)</a>
```

- 4、代码实现

```
1 from urllib import request
2 import re
3 from useragents import ua_list
4 import time
5 import random
6
7 class FilmSkySpider(object):
```

```

8     def __init__(self):
9         # 一级页面url地址
10        self.url = 'https://www.dytt8.net/html/gndy/dyzz/list_23_{}.html'
11
12    # 获取html功能函数
13    def get_html(self,url):
14        headers = {
15            'User-Agent':random.choice(ua_list)
16        }
17        req = request.Request(url=url,headers=headers)
18        res = request.urlopen(req)
19        # 通过网站查看网页源码,查看网站charset='gb2312'
20        # 如果遇到解码错误,识别不了某些字符,则 ignore 忽略掉
21        html = res.read().decode('gb2312','ignore')
22
23        return html
24
25    # 正则解析功能函数
26    def re_func(self,re_bds,html):
27        pattern = re.compile(re_bds,re.S)
28        r_list = pattern.findall(html)
29
30        return r_list
31
32    # 获取数据函数 - html是一级页面响应内容
33    def parse_page(self,one_url):
34        html = self.get_html(one_url)
35        re_bds = r'<table width="100%" .*?<td width="5%" .*?<a href="(.*?)".*?</table>'
36        # one_page_list: ['/html/xxx','/html/xxx','/html/xxx']
37        one_page_list = self.re_func(re_bds,html)
38
39        for href in one_page_list:
40            two_url = 'https://www.dytt8.net' + href
41            self.parse_two_page(two_url)
42            # uniform: 浮点数,爬取1个电影信息后sleep
43            time.sleep(random.uniform(1, 3))
44
45
46    # 解析二级页面数据
47    def parse_two_page(self,two_url):
48        item = {}
49        html = self.get_html(two_url)
50        re_bds = r'<div class="title_all"><h1><font color=#07519a>(.*?)</font></h1></div>.*?<td
51        style="WORD-WRAP.*?>.*?>(.*?)</a>'
52        # two_page_list: [('名称1','ftp://xxxx.mkv')]
53        two_page_list = self.re_func(re_bds,html)
54
55        item['name'] = two_page_list[0][0].strip()
56        item['download'] = two_page_list[0][1].strip()
57
58        print(item)
59
60    def main(self):
61        for page in range(1,201):
62            one_url = self.url.format(page)
63            self.parse_page(one_url)

```

```

64         # uniform: 浮点数
65         time.sleep(random.uniform(1,3))
66
67     if __name__ == '__main__':
68         spider = FilmSkySpider()
69         spider.main()

```

■ 5、练习

把电影天堂数据存入MySQL数据库 - 增量爬取

```

1  # 思路
2  # 1、MySQL中新建表 urltab,存储所有爬取过的链接的指纹
3  # 2、在爬取之前,先判断该指纹是否爬取过,如果爬取过,则不再继续爬取

```

练习代码实现

```

1  # 建库建表
2  create database filmskydb charset utf8;
3  use filmskydb;
4  create table request_finger(
5  finger char(32)
6  )charset=utf8;
7  create table filmtab(
8  name varchar(200),
9  download varchar(500)
10 )charset=utf8;

```

```

1  from urllib import request
2  import re
3  from useragents import ua_list
4  import time
5  import random
6  import pymysql
7  from hashlib import md5
8
9
10 class FilmSkySpider(object):
11     def __init__(self):
12         # 一级页面url地址
13         self.url = 'https://www.dytt8.net/html/gndy/dyzz/list_23_{}.html'
14         self.db = pymysql.connect('192.168.153.151', 'tiger', '123456', 'filmskydb',
charset='utf8')
15         self.cursor = self.db.cursor()
16
17     # 获取html功能函数
18     def get_html(self, url):
19         headers = {
20             'User-Agent': random.choice(ua_list)
21         }
22         req = request.Request(url=url, headers=headers)
23         res = request.urlopen(req)
24         # 通过网站查看网页源码,查看网站charset='gb2312'
25         # 如果遇到解码错误,识别不了某些字符,则 ignore 忽略掉
26         html = res.read().decode('gb2312', 'ignore')

```

```

27
28         return html
29
30     # 正则解析功能函数
31     def re_func(self, re_bds, html):
32         pattern = re.compile(re_bds, re.S)
33         r_list = pattern.findall(html)
34
35         return r_list
36
37     # 获取数据函数
38     def parse_page(self, one_url):
39         html = self.get_html(one_url)
40         re_bds = r'<table width="100%".*?<td width="5%".*?<a href="(.*?)".*?uLink">.*?
</table>'
41         # one_page_list: ['/html/xxx', '/html/xxx', '/html/xxx']
42         one_page_list = self.re_func(re_bds, html)
43
44         for href in one_page_list:
45             two_url = 'https://www.dytt8.net' + href
46             # 生成指纹 - md5加密
47             s = md5()
48             s.update(two_url.encode())
49             two_url_md5 = s.hexdigest()
50             # 判断链接是否需要抓取
51             if self.is_go_on(two_url_md5):
52                 self.parse_two_page(two_url)
53                 # 爬取完成此链接后将指纹放到数据库表中
54                 ins = 'insert into request_finger values(%s)'
55                 self.cursor.execute(ins, [two_url_md5])
56                 self.db.commit()
57                 # uniform: 浮点数,爬取1个电影信息后sleep
58                 time.sleep(random.uniform(1, 3))
59
60
61     def is_go_on(self, two_url_md5):
62         # 爬取之前先到数据库中查询比对
63         sel = 'select finger from request_finger where finger=%s'
64         # 开始抓取之前,先来判断该链接之前是否抓取过
65         result = self.cursor.execute(sel, [two_url_md5])
66         if not result:
67             return True
68
69
70     # 解析二级页面数据
71     def parse_two_page(self, two_url):
72         item = {}
73         html = self.get_html(two_url)
74         re_bds = r'<div class="title_all"><h1><font color=#07519a>(.*?)</font></h1>
</div>.*?<td style="WORD-WRAP.*?>.*?>(.*?)</a>'
75         # two_page_list: [('名称1', 'ftp://xxxx.mkv')]
76         two_page_list = self.re_func(re_bds, html)
77
78         item['name'] = two_page_list[0][0].strip()
79         item['download'] = two_page_list[0][1].strip()
80
81         ins = 'insert into filmtab values(%s,%s)'

```

```

82         film_list = [
83             item['name'], item['download']
84         ]
85         self.cursor.execute(ins, film_list)
86         self.db.commit()
87         print(film_list)
88
89
90     def main(self):
91         for page in range(1, 201):
92             one_url = self.url.format(page)
93             self.parse_page(one_url)
94             # uniform: 浮点数
95             time.sleep(random.uniform(1, 3))
96
97
98 if __name__ == '__main__':
99     spider = FilmSkySpider()
100    spider.main()

```

今日作业

- 1 1、电影天堂数据,存入MySQL、MongoDB、CSV文件
- 2 2、百度图片抓取: 输入要抓取的图片内容,抓取首页的30张图片,保存到对应的文件夹,比如:
3 你想要谁的照片, 请输入: 赵丽颖
- 4 创建文件夹到指定目录: 赵丽颖 并把首页30张图片保存到此文件夹下
- 5 3、抓取链家二手房房源信息(房源名称、总价),把结果存入到MySQL数据库,MongoDB数据库,CSV文件
- 6 # 小区名、总价、单价