

# 王伟超

[wangweichao@tedu.cn](mailto:wangweichao@tedu.cn)

## DAY01

### 网络爬虫概述

#### • 定义

- 1 网络蜘蛛、网络机器人，抓取网络数据的程序
- 2 其实就是用Python程序模仿人点击浏览器并访问网站，而且模仿的越像越好，让web站点无法发现你不是人

#### • 爬取数据目的

- 1 1、公司项目测试数据
- 2 2、公司业务部门及其他部门所需数据
- 3 3、数据分析

#### • 企业获取数据方式

- 1 1、公司自有数据
- 2 2、第三方数据平台购买（数据堂、贵阳大数据交易所）
- 3 3、爬虫爬取数据

#### • Python做爬虫优势

- 1 1、Python：请求模块、解析模块丰富成熟，强大的Scrapy网络爬虫框架
- 2 2、PHP：对多线程、异步支持不太好
- 3 3、JAVA：代码笨重，代码量大
- 4 4、C/C++：虽然效率高，但是代码成型慢

#### • 爬虫分类

- 1 1、通用网络爬虫（搜索引擎使用，遵守robots协议）  
robots协议：网站通过robots协议告诉搜索引擎哪些页面可以抓取，哪些页面不能抓取，通用网络爬虫需要遵守robots协议（君子协议）  
<https://www.taobao.com/robots.txt>
- 2 2、聚焦网络爬虫：自己写的爬虫程序

#### • 爬虫爬取数据步骤

- 1 1、确定需要爬取的URL地址
- 2 2、由请求模块向URL地址发出请求,并得到网站的响应
- 3 3、从响应内容中提取所需数据
- 4     1、所需数据,保存
- 5     2、页面中有其他需要继续跟进的URL地址,继续第2步去发请求,如此循环

## 爬虫请求模块一

### 模块名及导入

- 1 1、模块名: urllib.request
- 2 2、导入方式:
- 3     1、import urllib.request
- 4     2、from urllib import request

### 常用方法详解

#### urllib.request.urlopen

- 作用

向网站发起请求并获取响应对象

- 参数

- 1 1、url: 需要爬取的URL地址
- 2 2、timeout: 设置等待超时时间,指定时间内未得到响应抛出超时异常

- 第一个爬虫程序 - 01\_urlopen.py

打开浏览器,输入百度地址(<http://www.baidu.com/>),得到百度的响应

```
1 import urllib.request
2
3 # urlopen() : 向URL发请求,返回响应对象
4 response=urllib.request.urlopen('http://www.baidu.com/')
5 # 提取响应内容
6 html = response.read().decode('utf-8')
7 # 打印响应内容
8 print(html)
```

- 响应对象 (response) 方法

- 1 1、bytes = response.read() # read()得到结果为 bytes 数据类型
- 2 2、string = response.read().decode() # decode() 转为 string 数据类型
- 3 3、url = response.geturl() # 返回实际数据的URL地址
- 4 4、code = response.getcode() # 返回HTTP响应码
- 5 # 补充
- 6 5、string.encode() # bytes -> string
- 7 6、bytes.decode() # string -> bytes

思考：网站如何来判定是人类正常访问还是爬虫程序访问？？？

```
1 # 向测试网站: http://httpbin.org/get 发请求,查看自己请求头
2 # 代码如下
3 import urllib.request
4 response=urllib.request.urlopen('http://httpbin.org/get')
5 html = response.read().decode()
6 print(html)
7
8 # html中的请求头headers如下
9 "headers": {
10     "Accept-Encoding": "identity",
11     "Host": "httpbin.org",
12     "User-Agent": "Python-urllib/3.6"
13 },
14 发现请求头中User-Agent竟然是:Python-urllib/3.6!!!!!!!!!!!!!!!!!!!!!!
15 我们需要重构User-Agent,发请求时带着User-Agent过去,但是 urllib.urlopen()方法不支持重构User-Agent,那我们
    怎么办? 请看下面的方法!!!
```

## urllib.request.Request

- 作用

创建请求对象(包装请求, 重构User-Agent, 使程序更像正常人类请求)

- 参数

```
1 1、url: 请求的URL地址
2 2、headers: 添加请求头 (爬虫和反爬虫斗争的第一步)
```

- 使用流程

```
1 1、构造请求对象(重构User-Agent)
2   req = urllib.request.Request(
3       url = 'http://httpbin.org/get'
4       headers={'User-Agent':'Mozilla/5.0'}
5   )
6 2、发请求获取响应对象(urlopen)
7   res = urllib.request.urlopen(req)
8 3、获取响应对象内容
9   html = res.read().decode('utf-8')
```

- 示例 - 02\_Request.py

向测试网站 (<http://httpbin.org/get>) 发起请求, 构造请求头并从响应中确认请求头信息

```
1 from urllib import request
2
3 # 定义常用变量: URL 、 headers
4 url = 'http://httpbin.org/get'
5 headers = {
6     'User-Agent':'Opera/9.80 (Windows NT 6.1; U; zh-cn) Presto/2.9.168 Version/11.50'
```

```

7 | }
8 | # 1. 创建请求对象 - 包装,并没有真正发请求
9 | req = request.Request(url=url,headers=headers)
10 | # 2. 获取响应对象
11 | res = request.urlopen(req)
12 | # 3. 提取响应内容
13 | html = res.read().decode('utf-8')
14 |
15 | print(html)

```

## URL地址编码模块

### 模块名及导入

- 模块

```

1 | # 模块名
2 | urllib.parse
3 |
4 | # 导入
5 | import urllib.parse
6 | from urllib import parse

```

- 作用

给URL地址中查询参数进行编码

```

1 | 编码前: https://www.baidu.com/s?wd=美女
2 | 编码后: https://www.baidu.com/s?wd=%E7%BE%8E%E5%A5%B3

```

### 常用方法

#### urllib.parse.urlencode({dict})

- URL地址中一个查询参数

```

1 | # 查询参数: {'wd' : '美女'}
2 | # urlencode编码后: 'wd=%e7%be%8e%e5%a5%b3'
3 |
4 | # 示例代码
5 | query_string = {'wd' : '美女'}
6 | result = urllib.parse.urlencode(query_string)
7 | # result: 'wd=%e7%be%8e%e5%a5%b3'

```

- URL地址中多个查询参数

```

1 from urllib import parse
2 params = {
3     'wd' : '美女',
4     'pn' : '50'
5 }
6 params = parse.urlencode(query_string_dict)
7 url = 'http://www.baidu.com/s?{}'.format(params)
8 print(url)

```

- 拼接URL地址的3种方式

```

1 # 1、字符串相加
2 baseurl = 'http://www.baidu.com/s?'
3 params = 'wd=%E7XXXX&pn=20'
4 url = baseurl + params
5
6 # 2、字符串格式化（占位符）
7 params = 'wd=%E7XXXX&pn=20'
8 url = 'http://www.baidu.com/s?%s'% params
9
10 # 3、format()方法
11 url = 'http://www.baidu.com/s?{'
12 params = 'wd=%E7XXXX&pn=20'
13 url = url.format(params)

```

- 练习 在百度中输入要搜索的内容，把响应内容保存到本地文件

```

1 请输入搜索内容：赵丽颖
2 # 最终保存到本地文件 - 赵丽颖.html

```

## 代码实现 - 03\_parse\_baidu.py

```

1 from urllib import request
2 from urllib import parse
3
4 # 拼接URL地址
5 def get_url(word):
6     url = 'http://www.baidu.com/s?{'
7     # params: wd=%E7XXXX
8     params = parse.urlencode({'wd':word})
9     url = url.format(params)
10
11     return url
12
13
14 # 发请求,保存本地文件
15 def request_url(url,filename):
16     headers = {'User-Agent':'Mozilla/5.0'}
17     # 请求对象 + 响应对象 + 提取内容
18     req = request.Request(url=url,headers=headers)
19     res = request.urlopen(req)

```

```

20     html = res.read().decode('utf-8')
21     # 保存数据
22     with open(filename, 'w', encoding='utf-8') as f:
23         f.write(html)
24
25     # 主程序入口
26     if __name__ == '__main__':
27         word = input('请输入搜索内容:')
28         url = get_url(word)
29         filename = word + '.html'
30         request_url(url, filename)

```

## quote(string)编码

- 示例1

```

1  from urllib import parse
2
3  string = '美女'
4  print(parse.quote(string))
5  # 结果: %E7%BE%8E%E5%A5%B3

```

改写之前urlencode()代码，使用quote()方法实现

```

1  from urllib import parse
2
3  url = 'http://www.baidu.com/s?wd={}'
4  word = input('请输入要搜索的内容:')
5  query_string = parse.quote(word)
6  print(url.format(query_string))

```

## unquote(string)解码

- 示例

```

1  from urllib import parse
2
3  string = '%E7%BE%8E%E5%A5%B3'
4  result = parse.unquote(string)
5  print(result)

```

## 总结

```

1  # 1、urllib.request
2  req = urllib.request.Request(url, headers)
3  res = urllib.request.urlopen(req)
4  html = res.read().decode('utf-8')
5
6  # 2、响应对象res方法
7  res.read()
8  res.getcode()

```

```
9 res.geturl()
10
11 # 3、urllib.parse
12 urllib.parse.urlencode({})
13 urllib.parse.quote(string)
14 urllib.parse.unquote(string)
```

## 百度贴吧数据抓取案例

### 要求

```
1 1、输入贴吧名称:赵丽颖吧
2 2、输入起始页:1
3 3、输入终止页:3
4 4、保存到本地文件
5 赵丽颖吧-第1页.html、赵丽颖吧-第2页.html ...
```

### 实现步骤

- 1、查看是否为静态页面

```
1 右键 - 查看网页源代码 - 搜索数据关键字
```

- 2、找URL规律

```
1 第1页:http://tieba.baidu.com/f?kw=? ? &pn=0
2 第2页:http://tieba.baidu.com/f?kw=? ? &pn=50
3 第n页:pn=(n-1)*50
```

- 3、获取网页内容
- 4、提取所需数据
- 5、保存(本地文件、数据库)

代码实现 - 04\_tieba\_spider.py

```
1 from urllib import request,parse
2 import time
3 import random
4 from useragents import ua_list
5
6 class BaiduSpider(object):
7     def __init__(self):
8         self.url = 'http://tieba.baidu.com/f?kw={}&pn={}'
9
10    # 获取
11    def get_html(self,url):
12        headers = {
13            'User-Agent': random.choice(ua_list)
14        }
```

```

15     req = request.Request(url=url,headers=headers)
16     res = request.urlopen(req)
17     html = res.read().decode('utf-8')
18
19     return html
20
21     # 解析
22     def parse_html(self):
23         pass
24
25     # 保存
26     def write_html(self, filename, html):
27         with open(filename, 'w', encoding='utf-8') as f:
28             f.write(html)
29
30     # 入口函数
31     def run(self):
32         name = input('请输入贴吧名:')
33         begin = int(input('请输入起始页:'))
34         end = int(input('请输入终止页:'))
35
36         # 查询参数编码
37
38         params = parse.quote(name)
39         for page in range(begin, end+1):
40             # URL拼接
41             pn = (page-1)*50
42             pn = (page-1)*50
43             url = self.url.format(params, pn)
44
45             # 调用类内函数进行页面抓取+保存
46             filename = '{}-第{}页.html'.format(name, page)
47             html = self.get_html(url)
48             self.write_html(filename, html)
49             # 控制爬取速度
50             time.sleep(random.randint(1, 3))
51             print('第%d页爬取完成' % page)
52
53     if __name__ == '__main__':
54         start = time.time()
55         spider = BaiduSpider()
56         spider.run()
57         end = time.time()
58         print('执行时间: %.2f' % (end-start))

```

## 正则解析模块

### re模块使用流程

- 方法一

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



- 方法二

```
1 # 1、创建正则编译对象
2 pattern = re.compile(r'正则表达式',re.S)
3 r_list = pattern.findall(html)
```

## 正则表达式元字符

元字符	含义
.	任意一个字符（不包括\n）
\d	一个数字
\s	空白字符
\S	非空白字符
[]	包含[]内容
*	出现0次或多次
+	出现1次或多次

思考：请写出匹配任意一个字符的正则表达式？

```
1 import re
2 # 方法一
3 pattern = re.compile('.',re.S)
4 # 方法二
5 pattern = re.compile('[\s\S]')
```

## 贪婪匹配和非贪婪匹配

- 贪婪匹配（默认）

```
1 1、在整个表达式匹配成功的前提下,尽可能多的匹配 *
```

```
2 2、表示方式：.*
```

- 非贪婪匹配

```
1 1、在整个表达式匹配成功的前提下,尽可能少的匹配 *
```

```
2 2、表示方式：.*?
```

### 示例代码 - 05\_re\_greed.py

```
1 import re
2
3 html = '''
4 <div><p>九霄龙吟惊天变</p></div>
```

```

5 <div><p>风云际汇潜水游</p></div>
6 '''
7 # 贪婪匹配 : .*
8 pattern = re.compile('<div><p>.*</p></div>', re.S)
9 r_list = pattern.findall(html)
10 print(r_list)
11
12 # 非贪婪匹配 : .*?
13 pattern = re.compile('<div><p>.*?</p></div>', re.S)
14 r_list = pattern.findall(html)
15 print(r_list)

```

## 正则表达式分组

- 作用

在完整的模式中定义子模式，将每个圆括号中子模式匹配出来的结果提取出来

- 示例

```

1 import re
2
3 s = 'A B C D'
4 p1 = re.compile('\w+\s+\w+')
5 print(p1.findall(s))
6 # 结果: ['A B', 'C D']
7
8 p2 = re.compile('(\w+)\s+\w+')
9 print(p2.findall(s))
10 # 结果: ['A', 'C']
11
12 p3 = re.compile('(\w+)\s+(\w+)')
13 print(p3.findall(s))
14 # 结果: [('A', 'B'), ('C', 'D')]

```

- 分组总结

```

1 1、在网页中,想要什么内容,就加()
2 2、先按整体正则匹配,然后再提取分组()中的内容
3 如果有2个及以上分组(),则结果中以元组形式显示 [('小区1', '500万'), ('小区2', '600万'), ()]

```

- 练习

页面结构如下:

```

1 # <div class="animal">.*?title="(.*?)".*?
2 <div class="animal">
3     <p class="name">
4         <a title="Tiger"></a>
5     </p>
6     <p class="content">
7         Two tigers two tigers run fast

```

```

8     </p>
9 </div>
10
11 <div class="animal">
12     <p class="name">
13         <a title="Rabbit"></a>
14     </p>
15
16     <p class="content">
17         Small white rabbit white and white
18     </p>
19 </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 - 06\_re\_exercise.py

```

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')
```

## 今日作业

1、把百度贴吧案例重写一遍,不要参照课上代码 2、爬取猫眼电影信息：猫眼电影-榜单-top100榜

```
1 | 第1步完成:
2 |   猫眼电影-第1页.html
3 |   猫眼电影-第2页.html
4 |   ... ..
5 |
6 | 第2步完成:
7 |   1、提取数据：电影名称、主演、上映时间
8 |   2、先打印输出,然后再写入到本地文件
```

3、复习任务

```
1 | pymysql、MySQL基本命令
2 | MySQL：建库建表普通查询等
```