

Spider_Day02回顾

爬取网站思路

- 1 1、先确定是否为动态加载网站
- 2 2、找URL规律
- 3 3、正则表达式
- 4 4、定义程序框架，补全并测试代码

数据持久化 - csv

```
1 import csv
2 with open('xxx.csv','w') as f:
3     writer = csv.writer(f)
4     writer.writerow([])
5     writer.writerows([(,),(),(),())])
```

数据持久化 - MySQL

```
1 import pymysql
2
3 # __init__(self):
4     self.db = pymysql.connect('IP',... ...)
5     self.cursor = self.db.cursor()
6 # write_data(self):
7     self.cursor.execute('sql',[data1])
8     self.cursor.executemany('sql',[(data1),(data2),(data3)])
9     self.db.commit()
10 # main(self):
11     self.cursor.close()
12     self.db.close()
```

数据持久化 - MongoDB

```
1 import pymongo
2
```

```

3  # __init__(self):
4      self.conn = pymongo.MongoClient('IP',27017)
5      self.db = self.conn['db_name']
6      self.myset = self.db['set_name']
7  # write_data(self):
8      self.myset.insert_one(dict)
9
10 # MongoDB - Command
11 >show dbs
12 >use db_name
13 >show collections
14 >db.collection_name.find().pretty()
15 >db.collection_name.count()
16 >db.collection_name.drop()
17 >db.dropDatabase()

```

多级页面数据抓取

```

1  # 整体思路
2  1、爬取一级页面,提取 所需数据+链接,继续跟进
3  2、爬取二级页面,提取 所需数据+链接,继续跟进
4  3、... ...
5  # 代码实现思路
6  1、所有数据最终都会在一级页面解析函数中拿到
7  2、避免重复代码 - 请求、解析需定义函数

```

requests模块

■ get()

```

1  1、发请求并获取响应对象
2  2、res = requests.get(url,headers=headers)

```

■ 响应对象res属性

```

1  res.text : 字符串
2  res.content : bytes
3  res.encoding: 字符编码 res.encoding='utf-8'
4  res.status_code : HTTP响应码
5  res.url : 实际数据URL地址

```

■ 非结构化数据保存

```

1  with open('xxx.jpg','wb') as f:
2      f.write(res.content)

```

Chrome浏览器安装插件

■ 安装方法

```
1 # 在线安装
2 1、下载插件 - google访问助手
3 2、安装插件 - google访问助手: Chrome浏览器-设置-更多工具-扩展程序-开发者模式-拖拽(解压后的插件)
4 3、在线安装其他插件 - 打开google访问助手 - google应用商店 - 搜索插件 - 添加即可
5
6 # 离线安装
7 1、下载插件 - xxx.crx 重命名为 xxx.zip
8 2、输入地址: chrome://extensions/ 打开- 开发者模式
9 3、拖拽 插件(或者解压后文件夹) 到浏览器中
10 4、重启浏览器, 使插件生效
```

■ 爬虫常用插件

```
1 1、google-access-helper : 谷歌访问助手,可访问 谷歌应用商店
2 2、Xpath Helper: 轻松获取HTML元素的XPath路径
3 开启/关闭: Ctrl + Shift + x
4 3、JsonView: 格式化输出json格式数据
5 4、Proxy SwitchyOmega: Chrome浏览器中的代理管理扩展程序
```

Spider_Day03笔记

xpath解析

■ 定义

```
1 XPath即为XML路径语言,它是一种用来确定XML文档中某部分位置的语言,同样适用于HTML文档的检索
```

■ 示例HTML代码

```
1 <ul class="CarList">
2   <li class="bjd" id="car_001" href="http://www.bjd.com/">
3     <p class="name">布加迪</p>
4     <p class="model">威航</p>
5     <p class="price">2500万</p>
6     <p class="color">红色</p>
7   </li>
8
9   <li class="byd" id="car_002" href="http://www.byd.com/">
10    <p class="name">比亚迪</p>
11    <p class="model">秦</p>
12    <p class="price">15万</p>
13    <p class="color">白色</p>
14  </li>
```

```
15 | </ul>
```

■ 匹配演示

```
1 | 1、查找所有的li节点
2 | //li
3 | 2、获取所有汽车的名称：所有li节点下的子节点p的值（class属性值为name）
4 | //li/p[@class="name"]
5 | 3、找比亚迪车的信息：获取ul节点下第2个li节点的汽车信息
6 | //ul[@class="CarList"]/li[2]/p
7 | 4、获取所有汽车的链接：ul节点下所有li子节点的href属性的值
8 | //ul[@class="CarList"]/li/@href
9 |
10 | # 只要涉及到条件,加 []
11 | # 只要获取属性值,加 @
```

■ 选取节点

```
1 | 1、// ：从所有节点中查找（包括子节点和后代节点）
2 | 2、@ ：获取属性值
3 | # 使用场景1（属性值作为条件）
4 | //div[@class="movie-item-info"]
5 | # 使用场景2（直接获取属性值）
6 | //div[@class="movie-item-info"]/a/img/@src
```

■ 匹配多路径（或）

```
1 | xpath表达式1 | xpath表达式2 | xpath表达式3
```

■ 常用函数

```
1 | 1、contains()：匹配属性值中包含某些字符串节点
2 | # 查找id属性值中包含字符串 "car_" 的 li 节点
3 | //li[contains(@id,"car_")]
4 | 2、text()：获取节点的文本内容
5 | # 查找所有汽车的价格
6 | //ul[@class="CarList"]/li/p[@class="price"]/text()
```

lxml解析库

■ 安装

```
1 | sudo pip3 install lxml
```

■ 使用流程

```

1 1、导模块
2   from lxml import etree
3 2、创建解析对象
4   parse_html = etree.HTML(html)
5 3、解析对象调用xpath
6   r_list = parse_html.xpath('xpath表达式')

```

■ html样本

```

1 <div class="wrapper">
2   <a href="/" id="channel">新浪社会</a>
3   <ul id="nav">
4     <li><a href="http://domestic.sina.com/" title="国内">国内</a></li>
5     <li><a href="http://world.sina.com/" title="国际">国际</a></li>
6     <li><a href="http://mil.sina.com/" title="军事">军事</a></li>
7     <li><a href="http://photo.sina.com/" title="图片">图片</a></li>
8     <li><a href="http://society.sina.com/" title="社会">社会</a></li>
9     <li><a href="http://ent.sina.com/" title="娱乐">娱乐</a></li>
10    <li><a href="http://tech.sina.com/" title="科技">科技</a></li>
11    <li><a href="http://sports.sina.com/" title="体育">体育</a></li>
12    <li><a href="http://finance.sina.com/" title="财经">财经</a></li>
13    <li><a href="http://auto.sina.com/" title="汽车">汽车</a></li>
14  </ul>
15 </div>

```

■ 示例+练习

```

1 from lxml import etree
2
3 html = '''
4 <div class="wrapper">
5   <a href="/" id="channel">新浪社会</a>
6   <ul id="nav">
7     <li><a href="http://domestic.sina.com/" title="国内">国内</a></li>
8     <li><a href="http://world.sina.com/" title="国际">国际</a></li>
9     <li><a href="http://mil.sina.com/" title="军事">军事</a></li>
10    <li><a href="http://photo.sina.com/" title="图片">图片</a></li>
11    <li><a href="http://society.sina.com/" title="社会">社会</a></li>
12    <li><a href="http://ent.sina.com/" title="娱乐">娱乐</a></li>
13    <li><a href="http://tech.sina.com/" title="科技">科技</a></li>
14    <li><a href="http://sports.sina.com/" title="体育">体育</a></li>
15    <li><a href="http://finance.sina.com/" title="财经">财经</a></li>
16    <li><a href="http://auto.sina.com/" title="汽车">汽车</a></li>
17  </ul>
18 </div>'''
19 # 创建解析对象
20 parse_html = etree.HTML(html)
21 # 调用xpath返回结果,text()为文本内容
22 r_list = parse_html.xpath('//a/text()')
23 #print(rList)
24
25 # 提取所有的href的属性值
26 r2 = parse_html.xpath('//a/@href')
27 #print(r2)

```

```

28
29 # 提取所有href的值,不包括 /
30 r3 = parse_html.xpath('//ul[@id="nav"]/li/a/@href')
31 #print(r3)
32
33 # 获取 图片、军事、...,不包括新浪社会
34 r4 = parse_html.xpath('//ul[@id="nav"]/li/a/text()')
35 for r in r4:
36     print(r)

```

xpath最常使用方法

```

1 1、先匹配节点对象列表
2   # r_list: ['节点对象1','节点对象2']
3   r_list = parse_html.xpath('基准xpath表达式')
4 2、遍历每个节点对象,利用节点对象继续调用 xpath
5   for r in r_list:
6       name = r.xpath('./xxxxxx')
7       star = r.xpath('./xxxxxx')
8       time = r.xpath('./xxxxxx')

```

链家二手房案例 (xpath)

■ 实现步骤

1. 确定是否为静态

```

1 打开二手房页面 -> 查看网页源码 -> 搜索关键字

```

2. xpath表达式

```

1 1、基准xpath表达式(匹配每个房源信息节点列表)
2   //ul[@class="sellListContent"]/li[@class="clear LOGVIEWDATA LOGCLICKDATA"]
3
4 2、依次遍历后每个房源信息xpath表达式
5   * 名称: './a[@data-el="region"]/text()'
6
7   # 户型+面积+方位+是否精装
8   info_list = './div[@class="houseInfo"]/text()' [0].strip().split('|')
9   * 户型: info_list[1]
10  * 面积: info_list[2]
11  * 方位: info_list[3]
12  * 精装: info_list[4]
13
14
15  * 楼层: './div[@class="positionInfo"]/text()'
16  * 区域: './div[@class="positionInfo"]/a/text()'
17  * 总价: './div[@class="totalPrice"]/span/text()'
18  * 单价: './div[@class="unitPrice"]/span/text()'

```

3. 代码实现

```

1 import requests
2 from lxml import etree
3 import time
4 import random
5
6 class LianjiaSpider(object):
7     def __init__(self):
8         self.url = 'https://bj.lianjia.com/ershoufang/pg{}/'
9         self.headers = {'User-Agent' : 'Mozilla/5.0'}
10
11     def get_page(self,url):
12         try:
13             # 设定超时时间,超时后抛出异常,被except捕捉,继续执行此函数再次请求
14             res = requests.get(url,headers=self.headers,timeout=5)
15             res.encoding = 'utf-8'
16             html = res.text
17             self.parse_page(html)
18         except Exception as e:
19             self.get_page(url)
20
21     def parse_page(self,html):
22         parse_html = etree.HTML(html)
23         # 基准xpath,匹配每个房源信息的节点对象
24         li_list = parse_html.xpath('//ul[@class="sellListContent"]/li[@class="clear LOGVIEWDATA LOGCLICKDATA"]')
25         # 定义空字典,用来存储抓取的最终数据
26         house_dict = {}
27         # 遍历依次匹配每个房源信息,获取所有所需数据
28         for li in li_list:
29             # 房源名称
30             name_list = li.xpath('.//a[@data-el="region"]/text()')
31             house_dict['house_name'] = [ name_list[0] if name_list else None ][0]
32
33             # 列表: 户型+面积+方位+是否精装
34             info_list = li.xpath('.//div[@class="houseInfo"]/text()')
35             house_info = [ info_list[0].strip().split('|') if info_list else None ][0]
36             if house_info:
37                 # 户型
38                 house_dict['house_model'] = house_info[1]
39                 # 面积
40                 house_dict['area'] = house_info[2]
41                 # 方位
42                 house_dict['direction'] = house_info[3]
43                 # 是否精装
44                 house_dict['hardcover'] = house_info[4]
45                 #####
46                 # 楼层
47                 floor_list = li.xpath('.//div[@class="positionInfo"]/text()')
48                 house_dict['floor'] = [ floor_list[0].strip()[:-2] if floor_list else None ][0]
49                 # 区域
50                 address_list = li.xpath('.//div[@class="positionInfo"]/a/text()')
51                 house_dict['address'] = [ address_list[0].strip() if address_list else None ][0]
52                 # 总价
53                 total_list = li.xpath('.//div[@class="totalPrice"]/span/text()')
54                 house_dict['total_price'] = [ total_list[0].strip() if total_list else None ][0]
55                 # 单价

```

```

56     unit_list = li.xpath('..//div[@class="unitPrice"]/span/text()')
57     house_dict['unit_price'] = [ unit_list[0].strip() if unit_list else None ][0]
58
59     print(house_dict)
60
61     def main(self):
62         for pg in range(1,11):
63             url = self.url.format(str(pg))
64             self.get_page(url)
65             print('第%d页爬取成功' % pg)
66             time.sleep(random.randint(1,3))
67
68 if __name__ == '__main__':
69     start = time.time()
70     spider = LianjiaSpider()
71     spider.main()
72     end = time.time()
73     print('执行时间:%.2f' % (end-start))

```

练习: 将猫眼电影用xpath去实现

■ 1、xpath表达式

```

1  1、基准xpath: 匹配所有电影信息的节点对象列表
2      //dl[@class="board-wrapper"]/dd
3
4  2、遍历对象列表, 依次获取每个电影信息
5      for dd in dd_list:
6          电影名称 : dd.xpath('./a/@title')[0].strip()
7          电影主演 : dd.xpath('..//p[@class="star"]/text()')[0].strip()
8          上映时间 : dd.xpath('..//p[@class="releasetime"]/text()')[0].strip()

```

■ 2、代码实现

```

1  # 1、将urllib库改为requests模块实现
2  # 2、改写parse_page()方法
3  def parse_page(self,html):
4      # 创建解析对象
5      parse_html = etree.HTML(html)
6      # 基准xpath节点对象列表
7      dd_list = parse_html.xpath('//dl[@class="board-wrapper"]/dd')
8      movie_dict = {}
9      # 依次遍历每个节点对象,提取数据
10     for dd in dd_list:
11         movie_dict['name'] = dd.xpath('./a/@title')[0].strip()
12         movie_dict['star'] = dd.xpath('..//p[@class="star"]/text()')[0].strip()
13         movie_dict['time'] = dd.xpath('..//p[@class="releasetime"]/text()')[0].strip()
14
15     print(movie_dict)

```

3、完整代码实现

```

1  import requests
2  from lxml import etree

```



```

3 import time
4 import random
5
6 class MaoyanSpider(object):
7     def __init__(self):
8         self.url = 'https://maoyan.com/board/4?offset={}'
9         self.headers = {'User-Agent' : 'Mozilla/5.0 (Windows NT 10.0; WOW64) AppleWebKit/537.36
(KHTML, like Gecko) Chrome/74.0.3729.169 Safari/537.36'}
10        # 添加计数(页数)
11        self.page = 1
12
13        # 获取页面
14        def get_page(self,url):
15            # random.choice一定要写在这里,每次请求都会随机选择
16            try:
17                res = requests.get(url,headers=self.headers)
18                res.encoding = 'utf-8'
19                html = res.text
20                self.parse_page(html)
21            except Exception as e:
22                print('Error')
23                self.get_page(url)
24
25        # 解析页面
26        def parse_page(self,html):
27            # 创建解析对象
28            parse_html = etree.HTML(html)
29            # 基准xpath节点对象列表
30            dd_list = parse_html.xpath('//dl[@class="board-wrapper"]/dd')
31            item = {}
32            # 依次遍历每个节点对象,提取数据
33            if dd_list:
34                for dd in dd_list:
35                    # 电影名称
36                    name_list = dd.xpath('.//p/a/@title')
37                    item['name'] = [ name_list[0].strip() if name_list else None ][0]
38                    # 电影主演
39                    star_list = dd.xpath('.//p[@class="star"]/text()')
40                    item['star'] = [ star_list[0].strip() if star_list else None ][0]
41                    # 上映时间
42                    time_list = dd.xpath('.//p[@class="releasetime"]/text()')
43                    item['time'] = [ time_list[0].strip() if time_list else None ]
44
45                    print(item)
46            else:
47                print('No Data')
48
49        # 主函数
50        def main(self):
51            for offset in range(0,31,10):
52                url = self.url.format(str(offset))
53                self.get_page(url)
54                print('第%d页完成' % self.page)
55                time.sleep(random.randint(1,3))
56                self.page += 1
57
58 if __name__ == '__main__':

```

```
59 spider = MaoyanSpider()
60 spider.main()
```

百度贴吧图片抓取

■ 目标

```
1  抓取指定贴吧所有图片
```

■ 思路

```
1  1、获取贴吧主页URL, 下一页, 找到不同页的URL规律
2  2、获取1页中所有帖子URL地址: [帖子链接1, 帖子链接2, ...]
3  3、对每个帖子链接发请求, 获取图片URL
4  4、向图片的URL发请求, 以wb方式写入本地文件
```

■ 实现步骤

1. 贴吧URL规律

```
1  http://tieba.baidu.com/f?kw=??&pn=50
```

2. xpath表达式

```
1  1、帖子链接xpath
2      //div[@class="t_con_clefix"]/div/div/div/a/@href
3
4  2、图片链接xpath
5      //div[@class="d_post_content j_d_post_content_clefix"]/img[@class="BDE_Image"]/@src
6
7  3、视频链接xpath
8      //div[@class="video_src_wrapper"]/embed/@data-video
9      # 注意: 此处视频链接前端对响应内容做了处理, 需要查看网页源代码来查看, 复制HTML代码在线格式化
```

3. 代码实现

```
1  import requests
2  from urllib import parse
3  from lxml import etree
4  import time
5  import random
6
7  class BaiduImgSpider(object):
8      def __init__(self):
9          self.url = 'http://tieba.baidu.com/f?{'
10         self.headers = {'User-Agent': 'Mozilla/4.0 (compatible; MSIE 8.0; Windows NT 6.1; WOW64; Trident/4.0; SLCC2; .NET CLR 2.0.50727; .NET CLR 3.5.30729; .NET CLR 3.0.30729; Media Center PC 6.0; .NET4.0C; InfoPath.3)'}
11
12         # 获取html函数
```

```

13 def get_html(self,url):
14     try:
15         res = requests.get(url=url,headers=self.headers)
16         res.encoding = 'utf-8'
17         html = res.text
18
19         return html
20     except Exception as e:
21         self.get_html(url)
22
23 # 解析html函数
24 def xpath_func(self,xpath_bds,html):
25     parse_html = etree.HTML(html)
26     r_list = parse_html.xpath(xpath_bds)
27
28     return r_list
29
30
31 # 一级页面:获取帖子链接,最终搞定所有图片下载
32 # 还记得吗? 多级页面抓取所有数据都在一级页面中搞定!!!
33 def get_tlink(self,url):
34     html = self.get_html(url)
35     xpath_bds = '//div[@class="t_con cleafix"]/div/div/div/a/@href'
36     # tlink_list: ['/p/23234','/p/9032323']
37     tlink_list = self.xpath_func(xpath_bds,html)
38     # 依次遍历每个帖子链接,搞定所有的图片下载
39     if tlink_list:
40         for tlink in tlink_list:
41             t_url = 'http://tieba.baidu.com' + tlink
42             # 提取图片链接并保存
43             self.get_image(t_url)
44             time.sleep(random.randint(1,3))
45     else:
46         print('No Data')
47
48 # 获取图片链接
49 def get_image(self,t_url):
50     html = self.get_html(t_url)
51     # 提取图片链接
52     xpath_bds = '//*[@class="d_post_content j_d_post_content clearfix"]/img/@src'
53     imglink_list = self.xpath_func(xpath_bds,html)
54
55     for imglink in imglink_list:
56         html_content = requests.get(imglink,headers=self.headers).content
57         filename = imglink[-10:]
58         with open(filename,'wb') as f:
59             f.write(html_content)
60             print('%s下载成功' % filename)
61
62 # 指定贴吧名称,起始页和终止页,爬取图片
63 def main(self):
64     name = input('请输入贴吧名:')
65     begin = int(input('请输入起始页:'))
66     end = int(input('请输入终止页:'))
67     for page in range(begin,end+1):
68         # 查询参数编码
69         params = {

```

```

70         'kw' : name,
71         'pn' : str( (page-1)*50 )
72     }
73     params = parse.urlencode(params)
74     url = self.url.format(params)
75
76     # 开始获取图片
77     self.get_tlink(url)
78
79 if __name__ == '__main__':
80     spider = BaiduImgSpider()
81     spider.main()

```

requests.get()参数

查询参数-params

■ 参数类型

1 字典,字典中键值对作为查询参数

■ 使用方法

```

1 1、 res = requests.get(url,params=params,headers=headers)
2 2、 特点:
3     * url为基准的url地址, 不包含查询参数
4     * 该方法会自动对params字典编码,然后和url拼接

```

■ 示例

```

1 import requests
2
3 baseurl = 'http://tieba.baidu.com/f?'
4 params = {
5     'kw' : '赵丽颖吧',
6     'pn' : '50'
7 }
8 headers = {'User-Agent' : 'Mozilla/4.0 (compatible; MSIE 8.0; Windows NT 6.1; WOW64;
9     Trident/4.0; SLCC2; .NET CLR 2.0.50727; .NET CLR 3.5.30729; .NET CLR 3.0.30729; Media Center
10    PC 6.0; .NET4.0C; InfoPath.3)'}
11 # 自动对params进行编码,然后自动和url进行拼接,去发请求
12 res = requests.get(baseurl,params=params,headers=headers)
13 res.encoding = 'utf-8'
14 print(res.text)

```

Web客户端验证 参数-auth

■ 作用及类型

- 1、针对于需要web客户端用户名密码认证的网站
- 2、auth = ('username', 'password')

■ 达内code课程方向案例

```
1 import requests
2 from lxml import etree
3 import os
4
5 class NoteSpider(object):
6     def __init__(self):
7         self.url = 'http://code.tarena.com.cn/AIDCode/aid1904/14-redis/'
8         self.headers = {'User-Agent': 'Mozilla/5.0'}
9         self.auth = ('tarenacode', 'code_2013')
10
11     # 获取
12     def get_html(self):
13         html = requests.get(url=self.url, auth=self.auth, headers=self.headers).text
14         return html
15
16     # 解析提取数据 + 把笔记压缩包下载完成
17     def parse_page(self):
18         html = self.get_html()
19         xpath_bds = '//a/@href'
20         parse_html = etree.HTML(html)
21         # r_list : ['../', 'day01', 'day02', 'redis_day01.zip']
22         r_list = parse_html.xpath(xpath_bds)
23         for r in r_list:
24             if r.endswith('.zip') or r.endswith('.rar'):
25                 print(r)
26
27 if __name__ == '__main__':
28     spider = NoteSpider()
29     spider.parse_page()
```

思考：爬取具体的笔记文件？

```
1 import requests
2 from lxml import etree
3 import os
4
5 class NoteSpider(object):
6     def __init__(self):
7         self.url = 'http://code.tarena.com.cn/AIDCode/aid1904/14-redis/'
8         self.headers = {'User-Agent': 'Mozilla/5.0'}
9         self.auth = ('tarenacode', 'code_2013')
10
11     # 获取
12     def get_html(self):
13         html = requests.get(url=self.url, auth=self.auth, headers=self.headers).text
14         return html
15
16     # 解析提取数据 + 把笔记压缩包下载完成
```

```

17     def parse_page(self):
18         html = self.get_html()
19         xpath_bds = '//a/@href'
20         parse_html = etree.HTML(html)
21         # r_list : ['../', 'day01', 'day02', 'redis_day01.zip']
22         r_list = parse_html.xpath(xpath_bds)
23         for r in r_list:
24             if r.endswith('.zip') or r.endswith('.rar'):
25                 file_url = self.url + r
26                 self.save_files(file_url, r)
27
28     def save_files(self, file_url, r):
29         html_content = requests.get(file_url, headers=self.headers, auth=self.auth).content
30         # 判断保存路径是否存在
31         directory = '/home/tarena/AID1904/redis/'
32         filename = directory + r
33         if not os.path.exists(directory):
34             os.makedirs(directory)
35
36         with open(filename, 'wb') as f:
37             f.write(html_content)
38             print(r, '下载成功')
39
40 if __name__ == '__main__':
41     spider = NoteSpider()
42     spider.parse_page()

```

SSL证书认证参数-verify

■ 适用网站及场景

- 1 1、适用网站：https类型网站但是没有经过 证书认证机构 认证的网站
- 2 2、适用场景：抛出 `SSLError` 异常则考虑使用此参数

■ 参数类型

```

1 1、verify=True(默认)    : 检查证书认证
2 2、verify=False (常用) : 忽略证书认证
3 # 示例
4 response = requests.get(
5     url=url,
6     params=params,
7     headers=headers,
8     verify=False
9 )

```

代理参数-proxies

▪ 定义

- 1、定义：代替你原来的IP地址去对接网络的IP地址。
- 2、作用：隐藏自身真实IP,避免被封。

▪ 普通代理

获取代理IP网站

- 1 西刺代理、快代理、全网代理、代理精灵、... ..

参数类型

- 1、语法结构

```
1 proxies = {
2     '协议': '协议://IP:端口号'
3 }
4 
```
- 2、示例

```
1 proxies = {
2     'http': 'http://IP:端口号',
3     'https': 'https://IP:端口号'
4 }
```

示例

1. 使用免费普通代理IP访问测试网站: <http://httpbin.org/get>

```
1 import requests
2
3 url = 'http://httpbin.org/get'
4 headers = {
5     'User-Agent': 'Mozilla/5.0'
6 }
7 # 定义代理,在代理IP网站中查找免费代理IP
8 proxies = {
9     'http': 'http://112.85.164.220:9999',
10    'https': 'https://112.85.164.220:9999'
11 }
12 html = requests.get(url,proxies=proxies,headers=headers,timeout=5).text
13 print(html)
```

2. 思考: 建立一个自己的代理IP池, 随时更新用来抓取网站数据

```
1 import requests
2 from lxml import etree
3 from fake_useragent import UserAgent
4
5
6 class GetProxyIP(object):
7     def __init__(self):
8         self.url = 'https://www.xicidaili.com/nn/'
9         self.test_url = 'http://www.baidu.com/'
10
```

```

11
12 # 生成随机的User-Agent
13 def get_random_header(self):
14     # 创建User-Agent对象
15     try:
16         ua = UserAgent()
17         headers = {'User-Agent':ua.random}
18         return headers
19     except Exception as e:
20         self.get_random_header()
21
22
23 # 从西刺代理网站上获取随机的代理IP
24 def get_ip_list(self):
25     headers = self.get_random_header()
26     # 访问西刺代理网站国内高匿代理, 找到所有的tr节点对象
27     html = requests.get(self.url, headers=headers).content.decode('utf-8')
28     parse_html = etree.HTML(html)
29     # 基准xpath, 匹配每个代理IP的节点对象列表
30     ipobj_list = parse_html.xpath('//tr')
31
32     # 从列表中第2个元素开始遍历, 因为第1个为: 字段名 (国家、IP、... ...)
33     with open('ip_port.txt','a') as f:
34         for ip in ipobj_list[1:]:
35             ip_info = ip.xpath('./td[2]/text()')[0]
36             port_info = ip.xpath('./td[3]/text()')[0]
37             ip_port = ip_info + ':' + port_info + '\n'
38
39             proxies = {
40                 'http': 'http://' + ip_info + ':' + port_info,
41                 'https': 'https://' + ip_info + ':' + port_info
42             }
43
44             if self.test_proxy_ip(proxies):
45                 f.write(ip_port)
46                 print(ip_port, 'success')
47             else:
48                 print(ip_port, 'failed')
49
50 # 测试抓取的代理IP是否可用
51 def test_proxy_ip(self,proxies):
52     try:
53         res =
54 requests.get(self.test_url,headers=self.get_random_header(),proxies=proxies,timeout=3)
55         if res.status_code == 200:
56             return True
57         except Exception as e:
58             return False
59
60 if __name__ == '__main__':
61     spider = GetProxyIP()
62     spider.get_ip_list()

```


今日作业

电影天堂 (xpath)

糗事百科 (xpath)

- 1、URL地址: <https://www.qiushibaike.com/text/>
- 2、目标 : 用户昵称、段子内容、好笑数量、评论数量