

Python爬虫实战一之使用Beautiful Soup抓...

开发需求：

- 1、爬取百度招聘下的招聘信息：岗位名称、工作地点、公司名称、薪资、发布时间
- 2、超链接招聘具体信息：职位描述
- 3、可以根据岗位名称包含主要字段进行分类存储
- 4、可以根据发布时间进行分类存储
- 5、可以根据新增水平进行分类存储

.....

开发思路

- 1、找到翻页的url规律

第一页：http://zhaopin.baidu.com/quanzhi?query=测试&sort_type=1&city=济南&detailmode=close&rn=20&pn=0

第二页：http://zhaopin.baidu.com/quanzhi?query=测试&sort_type=1&city=济南&detailmode=close&rn=20&pn=20

第三页：http://zhaopin.baidu.com/quanzhi?query=测试&sort_type=1&city=济南&detailmode=close&rn=20&pn=40

.....

- 2、找到超链接的url规律

<http://zhaopin.baidu.com'+href>

其中从招聘信息中获取href标签半个路径

- 3、如何获取需要的岗位名称、工作地点、公司名称、薪资、发布时间、职位描述等信息，并封装为一个字典对象中

- 4、生成excel文件，并把字典数据存储进去

```
1 import xlwt
```

```

2 import time,os
3 class StatisticsReport(object):
4     t=time.strftime('%Y%m%d%H%M%S',time.localtime(time.time()))
5     #设置单元格样式
6     def set_style(self,name,height,bold=False):
7         # 初始化样式
8         style = xlwt.XFStyle()
9         # 为样式创建字体
10        font = xlwt.Font()
11        font.name = name
12        font.bold = bold
13        font.color_index = 4
14        font.height = height
15        style.font = font
16        return style
17    def __createStatisticsReport__(self):
18        RunNo=self.t
19        reportname=RunNo+'.xls'
20        self.__setreportname__(reportname)
21        ReportFile=xlwt.Workbook()
22        #创建1个测试报告sheet页名称
23        ReportFile.add_sheet(u'测试类岗位',cell_overwrite_ok=True)
24        #-----写入按测试类岗位的信息表头
25        #岗位名称    地区    公司名称    薪资    发布时间    岗位信息
26        wr_tree = ReportFile.get_sheet(0)
27        row0=[u'岗位名称',u'地区',u'公司名称',u'薪资',u'发布时间',u'岗位信息']
28        #生成按测试类岗位的信息表头
29        for i in range(0,len(row0)):
30            wr_tree.write(0,i,row0[i],self.set_style('Times New
Roman',220,True))
31            reportpath=os.path.abspath("..")+'\\'
32            print reportpath+reportname
33            ReportFile.save(reportpath+reportname)
34        def __setreportname__(self,reportname):
35            self.reportname=reportname
36        def __getreportname__(self):
37            return self.reportname

```

2、获取招聘信息并存储excel文件

```

1 import urllib2
2 from bs4 import BeautifulSoup
3 import xlrd,os
4 from xlutils.copy import copy
5 from StatisticsReport import StatisticsReport
6 def GetJobinfo():
7     dict1={}
8     n=1
9     head = {}    #设置头
10

```

```

11     head['User-Agent'] = 'Mozilla/5.0 (Windows NT 10.0; WOW64; rv:52.0)
    Gecko/20100101 Firefox/52.0'
12     urlhead='http://zhaopin.baidu.com/quanzhi?query=测试&sort_type=1&city=
    济南&detailmode=close&rn=20&pn=';
13     for i in range(0,n):
14         #获取url全路径
15         get_url=urlhead+str(2*i*10);
16         #模拟浏览器, 定制http请求头
17         request=urllib2.Request(url=get_url,headers = head)
18         #模拟浏览器, 调用urlopen并传入Request对象, 将返回一个相关请求response对
    象
19         reponse=urllib2.urlopen(request)
20         #这个应答对象如同一个文件对象, 可以在Response中调用.read()读取这个文件信
    息
21         zhaop_html=reponse.read().decode('utf-8')
22         # UTF-8模式读取获取的页面信息标签和内容
23         zhaop_htmltables=BeautifulSoup(zhaop_html,'lxml');
24         #获取所有'a'标签以及内容
25         get_linka_list=zhaop_htmltables.find_all('a')
26         j=0
27         for alink in get_linka_list:
28             href=alink.get('href')
29             if href.find('/szzw?detailidx')===-1:
30                 pass
31             else:
32                 get_joburl='http://zhaopin.baidu.com'+href;
33                 jobname=alink.span.get_text()
34                 jobarea=alink.find_all('p')[0].get_text().split('|')[0]
35                 jobcompany=alink.find_all('p')[0].get_text().split('|')[1]
36                 jobpay=alink.find_all('p')[1].get_text()
37                 jobdate=alink.find_all('p')[2].get_text()
38                 get_jobrequest=urllib2.Request(url=get_joburl,headers =
    head)
39                 #模拟浏览器, 调用urlopen并传入Request对象, 将返回一个相关请求
    response对象
40                 get_jobreponse=urllib2.urlopen(get_jobrequest)
41                 #这个应答对象如同一个文件对象, 可以在Response中调用.read()读取这
    个文件信息
42                 get_job_html=get_jobreponse.read().decode('utf-8')
43                 # UTF-8模式读取获取的页面信息标签和内容
44                 zhaop_htmltables=BeautifulSoup(get_job_html,'lxml');
45                 job_tag=zhaop_htmltables.find_all(name='div',attrs=
    {'class':'abs'});
46                 job_decs=job_tag[0].find_all('p')
47                 job_decinfo=''
48                 for job_dec in job_decs:
49                     job_decinfo+=job_dec.get_text()+'\n'
50                 list1=
    [jobname,jobarea,jobcompany,jobpay,jobdate,job_decinfo]

```

```

51         j+=1
52         key='test'+str(j)
53         dict2=dict.fromkeys([key], list1)
54         dict1.update(dict2)
55     return dict1
56 def GenerateReport(report,job_dict):
57     #os.path.abspath("..") 用于返回上一级绝对路径
58     reportpath=os.path.abspath("..")+ '\\'
59     reportname=report.__getreportname__()
60     bk=xlrd.open_workbook(reportpath+reportname)
61     wb=copy(bk)
62     wa=wb.get_sheet(0)
63     for i in range(0,len(job_dict)):
64         for j in range(0,len(job_dict.values()[i])):
65             wa.write(i+1,j,job_dict.values()[i][j])
66     wb.save(reportpath+reportname)
67 if __name__ == '__main__':
68     # 方法调用
69     report=StatisticsReport()
70     report.__createStatisticsReport__()
71     job_dict=GetJobinfo()
72     GenerateReport(report,job_dict)

```

结果：

	A	B	C	D	E	F
1	岗位名称	地区	公司名称	薪资	发布时间	岗位信息
						职位类型：其它 发布时间：2017-07-24 有效日期：2018-07-24 岗位职责：无
2	济南良成月	济南市-历城	济南良成	面议	昨天发布	
3	LTE网优	济南-历城	山东闻远	4000-8000	今天发布	职位类型：其他技术职位发布时间：2017-07-25有效日期：2017-0
4	软件测试	济南-不限	软通动力	5000-8000	2017-07-23	职位类型：IT/互联网发布时间：2017-07-23有效日期：2017-09-02
5	自动化测试	济南	银客未来	12000-180	昨天发布	职位类型：技术发布时间：2017-07-24有效日期：2017-09-24
6	软件测试	济南	淄博智博	4000-8000	今天发布	职位类型：软件测试发布时间：2017-07-25有效日期：2017-08-24
7	LTE网优	济南-不限	山东闻远	3000-5000	昨天发布	职位类型：电子/通信发布时间：2017-07-24有效日期：2017-09-03
8	物流软件	济南-历下	山东海文	5000-8000	昨天发布	职位类型：采购/贸易/物流/交通运输发布时间：2017-07-24有效日期：2017-08-24
9	资深测试工	济南-历下	车满满	10000-180	今天发布	职位类型：自动化测试发布时间：2017-07-25有效日期：2017-08-24
10	测试工程师	济南-历城	山东思臣	4000-8000	今天发布	职位类型：测试经理发布时间：2017-07-25有效日期：2017-08-24
11	APP测试	济南-历下	北京东方	3000-6000	今天发布	职位类型：移动端测试发布时间：2017-07-25有效日期：2017-08-24
12	软件测试	济南-历下	北京锐融	5000-1000	今天发布	职位类型：软件测试发布时间：2017-07-25有效日期：2017-08-24
13	8279 软件	济南	深圳市金	7000-1200	今天发布	职位类型：技术发布时间：2017-07-25有效日期：2017-09-25
14	5千起聘	济南-历城	北京才秀	5000-8000	昨天发布	职位类型：计算机/互联网/通信发布时间：2017-07-24有效日期：2017-09-25
15	测试工程师	济南	深圳市金	8000-1200	今天发布	职位类型：技术发布时间：2017-07-25有效日期：2017-09-25
16	功能测试	济南-历下	山东泰商	3000-4000	今天发布	职位类型：功能测试发布时间：2017-07-25有效日期：2017-08-24
17	软件测试工	济南-不限	中软国际	5000-8000	昨天发布	职位类型：电子/通信发布时间：2017-07-24有效日期：2017-09-12
18	测试工程师	济南	北京辰森	6000-1200	今天发布	职位类型：技术发布时间：2017-07-25有效日期：2017-09-25
19	功能测试	济南-历下	山东泰商	3000-4000	今天发布	职位类型：功能测试发布时间：2017-07-25有效日期：2017-08-24
20	测试主管	济南-不限	山东世纪	5000-8000	2017-07-23	职位类型：IT/互联网发布时间：2017-07-23有效日期：2017-09-17
21	5K聘机械	济南-历下	山东海文	3000-5000	昨天发布	职位类型：生产制造/造纸印刷/服装服饰发布时间：2017-07-24有效日期：2017-08-24
22						
23						
24						
25						
26						
27						
28						
29						
30						