

实例：小型网站新闻爬虫词云图

19:13

一、目标分析

1选取目标

华东交大“交大要闻”新闻页面

<https://xw.ec.jtu.edu.cn/>

<https://xw.ecjtu.edu.cn/jdyw/684.htm>



2网页结构分析

(1) 页面目录

a,文章目录

第1页: <https://xw.ecjtu.edu.cn/jdyw.htm>

第2页: <https://xw.ecjtu.edu.cn/jdyw/684.htm>

第3页: <https://xw.ecjtu.edu.cn/jdyw/683.htm>

第4页: <https://xw.ecjtu.edu.cn/jdyw/682.htm>

第5页: <https://xw.ecjtu.edu.cn/jdyw/681.htm>

• • • • •

第685页: <https://xw.ecjtu.edu.cn/jdyw/1.htm>

b,文章页面

经过总结如右侧: <https://xw.ecjtu.edu.cn/info/1041/xxxxxx.htm>

(2) 同页结构

```
<div class="list fl">
<script language="javascript" src="/system/resource/js/centerOuting.js"></script><script language="javascript" src="/system/resource/js/ajax.js"></script><ul>
<li><a href="http://www.ecjtu.edu.cn/info/1053/30806.htm" target="_blank" title="【我为师生办实事】多举措推进良好校园人圈环境" >【我为师生办实事】多举措推进良好校园人圈环境</a></li>
<li><a href="http://www.ecjtu.edu.cn/info/1053/30805.htm" target="_blank" title="学校召开机械工程学院专业建设研讨会" >学校召开机械工程学院专业建设研讨会</a></li>
<li><a href="http://www.ecjtu.edu.cn/info/1053/30802.htm" target="_blank" title="上海市第十年度会计专业类十所高校大会" >上海市第十年度会计专业类十所高校大会</a></li>
<li><a href="http://www.ecjtu.edu.cn/info/1053/30800.htm" target="_blank" title="我校荣获2020年度国家社科基金项目" >我校荣获2020年度国家社科基金项目</a></li>
<li><a href="http://www.ecjtu.edu.cn/info/1053/30797.htm" target="_blank" title="我校在第十七届“挑战杯”全国大学生课外学术科技作品竞赛中获奖" >我校在第十七届“挑战杯”全国大学生课外学术科技作品竞赛中获奖</a></li>
<li><a href="http://www.ecjtu.edu.cn/info/1053/30793.htm" target="_blank" title="缅怀革命先烈 传承红色基因——我校开展党史红色主题教育活动" >缅怀革命先烈 传承红色基因——我校开展党史红色主题教育活动</a></li>
<li><a href="http://www.ecjtu.edu.cn/info/1053/30791.htm" target="_blank" title="学校召开党委理论学习中心组学习会" >学校召开党委理论学习中心组学习会</a></li>
<li><a href="http://www.ecjtu.edu.cn/info/1053/30788.htm" target="_blank" title="我校召开2020年中心工作会议暨2021年度工作会" >我校召开2020年中心工作会议暨2021年度工作会</a></li>
<li><a href="http://www.ecjtu.edu.cn/info/1053/30787.htm" target="_blank" title="学校召开50周年校庆总结大会" >学校召开50周年校庆总结大会</a></li>
<li><a href="http://www.ecjtu.edu.cn/info/1053/30785.htm" target="_blank" title="“互联网+”国际总决赛参赛的决赛暨颁奖典礼召开" >“互联网+”国际总决赛参赛的决赛暨颁奖典礼召开</a></li>
<li><a href="http://www.ecjtu.edu.cn/info/1053/30784.htm" target="_blank" title="我校通过连续两年地市级机械设计制造及其自动化、测绘工程专业工程教育认证" >我校通过连续两年地市级机械设计制造及其自动化、测绘工程专业工程教育认证</a></li>
<li><a href="http://www.ecjtu.edu.cn/info/1053/30783.htm" target="_blank" title="我校通过连续三年地市级机械设计制造及其自动化、测绘工程专业工程教育认证" >我校通过连续三年地市级机械设计制造及其自动化、测绘工程专业工程教育认证</a></li>
<li><a href="http://www.ecjtu.edu.cn/info/1053/30778.htm" target="_blank" title="我校通过连续四年地市级机械设计制造及其自动化、测绘工程专业工程教育认证" >我校通过连续四年地市级机械设计制造及其自动化、测绘工程专业工程教育认证</a></li>
<li><a href="http://www.ecjtu.edu.cn/info/1053/30776.htm" target="_blank" title="2021年度江西南自然科学基金项目申报结果公示" >2021年度江西南自然科学基金项目申报结果公示</a></li>
<li><a href="http://www.ecjtu.edu.cn/info/1053/30759.htm" target="_blank" title="京工党来交文大党支部委员会举行换届选举大会" >京工党来交文大党支部委员会举行换届选举大会</a></li>
</ul>
</div>
<link rel="stylesheet" Content-type="text/css" href="/system/resource/css/pagedown/sys.css"><div class="pub_sys_common pub_sys_normal pub_sys_style">
<div>
<div>
<div>
```

结构在list fl, <a>标签下

(3) 文章页面结构

```

<div class="content-con">
  <div id="sub_content"><div class="y_news_content">
    <p><span class="time">12</span></p><p>我校数学信息员队伍已实现一个，目前已评选出优秀数学信息员8名。</p><p></p><p>该制度的实施充分发挥了学生参与数学管理和自我管理的主力作用，加强了我校教师
    </div><p align="right">唐世编稿，<div id="y_vote"><div id="y_vote_id"></div>
  </div>

```

结构在content-con, <p>标签下

其他情况下存在渲染等非预期字符，故需要考虑白名单过滤汉字。

```
<div class="content-con">
  <div id="vsb_content"><div class="v_news_content">
    <span class="size12"><p><span style="font-size: 9pt; font-family: 宋体; ms-ascii-font-family: 'Times New Roman'; ms-ohang-font-family: 'Times
    </div><p align="right">责任编辑: admin</p></div><div id="div_vote_id"></div>
```

二，爬虫关键代码

1调用bs4库

```
from urllib.request import urlopen
```

```
from bs4 import BeautifulSoup
```

```
html = urlopen('url')
bs = BeautifulSoup(html.read(), 'html.parser')
#print(bs)
```

2爬取并截获文章内容

```
#爬取content-title的标题
nameList=bs.findAll('div',{'class':'content-title'})
h3_txt=[]
for i in nameList:
    h3_txt.append(str(i.findAll('h3'))))

#爬取文章内容(p)
nameList=bs.findAll('div',{'class':'v_news_content'})
p_txt=[]
for i in nameList:#由于tag的属性get_text需要针对每一个变量
    for j in i.findAll('span'):
        p_txt.append(str(j.get_text()))
```

3处理文章内容，并保存到列表

```
nameList=bs.findAll('h3')#标签爬取

strlist=[]
pattern = re.compile(r'[\u4e00-\u9fa5]+')#白名单过滤中文字符
for i in nameList:
    p = re.findall(pattern, str(i))
    if str(p)!='[]':#过滤为空的情况
        strlist.append(p)
    else:
        continue
print(strlist)
```

4爬取第一页url

```
import re
from urllib.request import urlopen
from bs4 import BeautifulSoup

html = urlopen('https://xw.ecjtu.edu.cn/jdyw.htm')#每页15篇文章内容
bs = BeautifulSoup(html.read(), 'html.parser')

#爬取第一页url
url_txt=[]
nameList=bs.findAll('div',{'class':'list fl'})
for i in nameList:
    for link in i.findAll('a'):
        if 'href' and 'target' in link.attrs:#根据基本格式抓取文章URL
            url_txt.append(link.attrs['href'])#保存到列表

print(url_txt)
```

三，词云关键代码

```
def img_grearte():
    mask = imread("xin.jpeg")
    with open("txt_save.txt", "r") as file:
        txt = file.read()
    word = WordCloud(background_color="white",
                      width=800,
                      height=800,
                      font_path='simhei.ttf',
                      mask=mask,
```



```
62 for i in range(5):#指定页数为5页
63     url="https://xw.ecjtu.edu.cn
64     f = open("C:\\Users\\qte\\Des
65     #print(url)
66     f.write(page_go(url))
67     f.close()
```

2词语脚本.py

修改禁用词:

```
62 for i in range(5):#指定页数为5页
63     url="https://xw.ecjtu.edu.cn
64     f = open("C:\\Users\\qte\\Des
65     #print(url)
66     f.write(page_go(url))
67     f.close()
```