

```

1 #201018_网络爬虫
2 #requests模块,这个模块可以爬取网页内容
3 import re
4 import requests
5 import xlwt,xlrd
6 url='https://www.xs4.cc/0_4/' #需要进行爬虫的网址
7 url2=re.findall('https://www.xs4.cc/(.*?)',url)[0] #提取书名所在网址的最后一段,后面要用
8 req=requests.get(url) #实例化一个requests请求,提取网页内容
9 req.encoding='gbk' #将编码转换为gbk,解决乱码问题
10 book_name=re.findall('<h1>(.*?)</h1>',req.text)[0] #获取书名
11 mulu=re.findall('.html">(.*?)</a></dd>',req.text,re.S)
12 # for i in range(9,len(mulu)):
13 #     print(mulu[i])
14 wangzhi=re.findall(f'<a href="{url2}/{url2}.html">',req.text,re.S)
15 dict1={}
16 for i in range(9,len(mulu)):
17     dict1[mulu[i]]=f'{url2}{wangzhi[i]}.html' #将目录和网址存放到字典中
18 # for i in range(9,len(wangzhi)):
19 #     print(f'{url2}{wangzhi[i]}.html')
20
21 # for k,v in dict1.items():
22 #     print(k,v)
23
24 excell=xlwt.workbook() #实例化一个excel
25 worksheet=excell.add_sheet(f'{book_name}') #新建一个sheet
26 worksheet.write(0,0,'目录') #行,列,内容
27 worksheet.write(0,1,'网址')
28 row=1
29 for k,v in dict1.items():
30     worksheet.write(row,0,k) #写入目录
31     worksheet.write(row,1,v) #写入网址
32     row+=1
33 excell.save(f'd:{book_name}.xls') #保存excel文件
34 #读取excel文件
35 data=xlrd.open_workbook(f'd:{book_name}.xls')
36 sheet1=data.sheets()[0] #读取文件的第一个sheet
37 # print(sheet1.nrows) #nrows返回 excel中的有效行数
38 for i in range(1,sheet1.nrows):
39     print(sheet1.cell_value(i,0),sheet1.cell_value(i,1)) #获取单元格内容
40
41 #上节课思考题,求某数的阶乘
42 def jiecheng(n):
43     if n==1:
44         return n
45     else:
46         return n*jiecheng(n-1)
47
48 #课后思考题,提取全书网正文

```