

```

1  #200823_网络爬虫
2  import re
3  import requests
4  import xlwt
5  import xlrd
6  web_url='https://www.xs4.cc/0_4/' #选择全书网的某一本书的网址
7  web_url2=re.findall('cc/(.*?)',web_url)[0] #取得书的网址的一部分的字符串
8  req=requests.get(web_url) #实例化一个requests请求
9  req.encoding='gbk' #编码设置为gbk,防止乱码
10 shuming=re.findall('<h1>(.*?)</h1>',req.text)[0] #取得书名
11 mulu=re.findall('html">(.*?)</a></dd>',req.text,re.S) #获取目录,生成一个列表
12 wangzhi=re.findall(f'<a href="{web_url2}/{(.*?)}.html',req.text,re.S) #获取书
    名的网址的一部分,也是一个列表
13 dict1={}
14 for i in range(9,len(mulu)): #从9开始循环,因为0-8是最新章节,之后才是第一章
15     dict1[mulu[i]]=f'{web_url}{wangzhi[i]}.html' #把目录和网址放到字典中
16
17 excel=xlwt.workbook() #实例化一个excel对象
18 worksheet=excel.add_sheet('全书网') #添加一个sheet
19 worksheet.write(0,0,'目录') #write(行,列,值)
20 worksheet.write(0,1,'地址')
21
22 row=1
23 for key,value in dict1.items():
24     worksheet.write(row,0,key) #将目录写入第0列
25     worksheet.write(row,1,value) #将地址写入第1列
26     row+=1
27 excel.save(f'd:{shuming}.xls') #保存excel文件
28 data=xlrd.open_workbook(f'd:{shuming}.xls') #读取excel
29 table=data.sheets()[0] #读取第一个sheet
30 for i in range(1,table.nrows): #table.nrows 有效行数
31     print(table.cell_value(i,0),table.cell_value(i,1)) #cell_value(行,列)获取
        单元格内容

```