

```

1  #201019_获取全书网正文
2  #上节课思考题 获取全书网正文
3  import re
4  import requests
5  import os
6  web_url='https://www.xs4.cc/0_4/' #选择要爬取的书的网址
7  web_url2=re.findall('cc/(.*?)',web_url)[0] #取得网址的一部分
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
12 wangzhi=re.findall(f'<a href="{web_url2}/{(.*?)}.html"',req.text) #获取网址的
    一部分
13 dict1={}
14 for i in range(9,len(mulu)):
15     dict1[mulu[i]]=f'{web_url}{wangzhi[i]}.html' #把目录和网址存入字典中
16 if os.path.exists(f'd://{shuming}'): #如果目录存在,则不做操作
17     pass
18 else:
19     os.mkdir(f'd://{shuming}') #如果目录不存在,则创建以书名命名的目录
20 count=0
21 for k,v in dict1.items():
22     if count>=3: #控制循环的次数
23         break
24     else:
25         zhengwen=requests.get(v)
26         zhengwen.encoding='gbk'
27         neirong=re.findall('<div id="content">(.*?)
    </div>',zhengwen.text,re.S)[0]
28         neirong=neirong.replace('&nbsp;','').replace('<br />','')
29         with open(f'd://{shuming}/{k}.txt','w+') as file1:
30             file1.write(neirong)
31         count += 1

```