

实验一：python 爬虫与 GUI 界面

PB18000169 陈晨曦

任务说明：

利用 python 实现网络爬虫，获取相应网站的所有图片，并利用图形化界面 GUI 来实现图片数据的获取与信息反馈。

实验细节：

get_jmp_strurl_list 函数：

函数的输入参数为字符串 str_url，即为网页的 URL；函数的返回为字符串数组，其中存放了相应网页的所有图片的 URL。

通过 python 自带的函数库 requests 和 BeautifulSoup 来实现 get_jmp_strurl_list 函数，在利用 select 函数获得网页的所有图片信息后，还需要将相应的图片信息利用 re 库来实现正则化处理，使得最后返回的 URL 符合规范。

源代码：

```
1 # -*- coding: utf-8 -*-
2 """
3 Created on Sun Apr 11 14:08:03 2021
4
5 @author: Lenovo
6 """
7 import requests
8 from bs4 import BeautifulSoup
9 import re
10
11 #从待爬网页获得图片内容
12 def get_jmp_strurl_list(str_url):
13     #str_url="https://www.ustc.edu.cn" #待爬网页的URL,这是一个字符串
14     URL=requests.get(str_url) #这是一个URL对象
15
16     #URL.text为URL对象中的html超文本文件内容
17     soup=BeautifulSoup(URL.text,'lxml') #转化成Unicode编码格式
18
19     #看看此时的strhtml,可见其文件结构满足html超文本文件的结构
20     #print(soup)
21
22     data=soup.select('img')
23     jmp_strurl_list=[]
24     for item in data:
25         result=str_url+"/"+item.get('src')
26         jmp_strurl_list.append(result)
27
28     #for jmp_strurl in jmp_strurl_list:
29         #print(jmp_strurl)
30     return jmp_strurl_list
```

save_img 函数：

函数的输入参数为 filename 和 jmp_strurl_list，其中 filename 代表要创建的文件夹的文件夹名（该文件夹用来存储下载的图片），jmp_strurl_list 是个字符串数组，存放要下载的图片的 URL。

为了实现该函数，首先要调用 os 库实现文件夹的创建和路径的设置，其次利用 requests 库来访问图片的 URL 并完成图片的下载。在实现该函数的过程中，还调用了 time

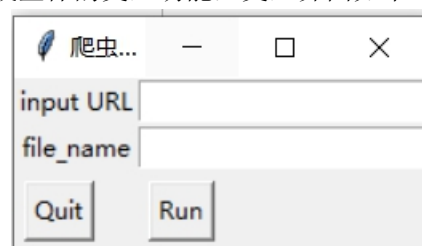
库来对无法访问 URL 进行处理，并在函数最后显示无法下载的图片的 URL。

源代码：

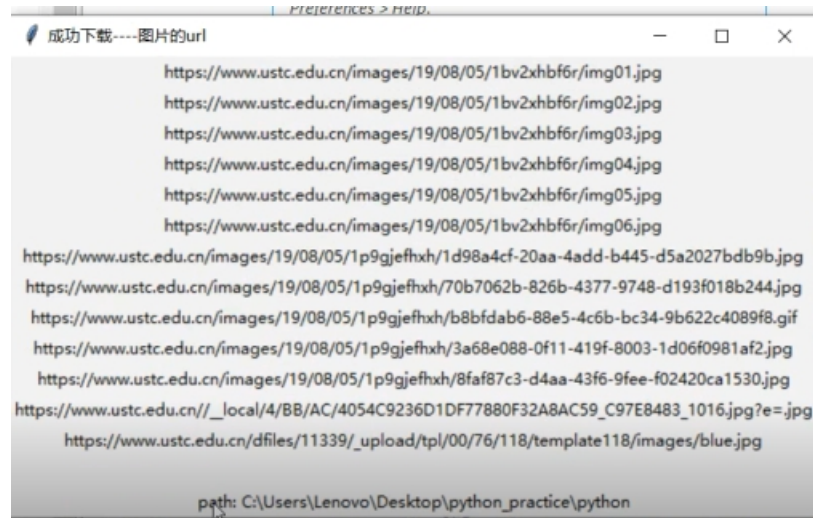
```
1 #-*- coding: utf-8 -*-
2 """
3 Created on Sun Apr 11 14:13:51 2021
4
5 @author: Lenovo
6 """
7 import requests
8 import time
9 import os
10
11 def save_img(file_name,jmp_strurl_list):
12     #创建保存下载图片的文件夹
13     count=0
14     mark=1
15     unable_jmp_strurl_list=[]
16     able_jmp_strurl_list=[]
17     if (os.path.exists("./"+file_name)):
18         print("the file_name already exists,please change the file_name")
19     else:
20         os.makedirs("./"+file_name)
21
22     #开始下载图片
23     for jmp_strurl in jmp_strurl_list:
24         print(" try downloading----- %s" %(jmp_strurl))
25         #通过时延语句判断图片是否可下载
26         try:
27             jmp_url=requests.get(jmp_strurl,timeout=10)
28             time.sleep(1)
29         except:
30             mark=0;
31
32     if(mark==1):
33         jmp_full_path="./"+file_name+"/"+str(count)+".jpg"
34         with open(jmp_full_path,'wb') as f:
35             f.write(jmp_url.content)
36         if not (jmp_strurl in able_jmp_strurl_list):
37             able_jmp_strurl_list.append(jmp_strurl)
38             count=count+1
39
40     elif(mark==0):
41         if not (jmp_strurl in unable_jmp_strurl_list):
42             unable_jmp_strurl_list.append(jmp_strurl)
43             mark=1
44
45     print()
46     print("successfully download %d pictures" %(count))
47     if(unable_jmp_strurl_list!=[]):
48         print("unable download jmp url list:")
49         for unable_jmp_strurl in unable_jmp_strurl_list:
50             print(" "+unable_jmp_strurl)
51
52     return able_jmp_strurl_list
53
```

图形化界面部分：

实现图形化交互界面，实现方法为调用 tkinter 库，创建窗口，利用库中的 Label 类，Button 类，Entry 类完成整体的交互功能，交互界面如下：



在 tkinter 的具体使用，通过窗口中再生成窗口的方法来实现图片数据的反馈处理，为了反馈下载成功的信息，我们生成新窗口 info_window：



在 info_window 中的最后一行，我们还通过 os.getcwd()命令返回了保存下载图片的文件夹的路径信息。

源代码：

```
1 # -*- coding: utf-8 -*-
2 """
3 Spyder Editor
4
5 This is a temporary script file.
6 """
7 import get_jmp
8 import download
9 import os
10
11 from tkinter import *
12
13 def GUI_func():
14     strurl=entry1.get()
15     file_name=entry2.get()
16     jmp_strurl_list=get_jmp_strurl_list(strurl)
17     able_jmp_strurl_list=save_img(file_name,jmp_strurl_list)
18
19     #显示成功下载的图片信息
20     info_window=Tk()
21     info_window.title("成功下载----图片的url")
22     i=0
23     for jmp_strurl in able_jmp_strurl_list:
24         Label(info_window,text=jmp_strurl).grid(row=i,column=0)
25         i=i+1
26
27     #窗口中空一行
28     Label(info_window,text="").grid(row=i,column=0)
29     i=i+1
30
31     #显示图片的路径
32     Label(info_window,text="path: "+os.getcwd()).grid(row=i,column=0)
33     i=i+1
34     pass
35
36 myWindow=Tk()
37 myWindow.title("爬虫----获取网页所有图片")
38 Label(myWindow,text="input URL").grid(row=0)
39 Label(myWindow,text="file_name").grid(row=1)
40
41 entry1=Entry(myWindow)
42 entry1.grid(row=0, column=1)
43 entry2=Entry(myWindow)
44 entry2.grid(row=1, column=1)
45
46 Button(myWindow, text='Quit', command=myWindow.quit).grid(row=2, column=0,sticky=W, padx=5, pady=5)
47 Button(myWindow, text='Run', command=GUI_func).grid(row=2, column=1, sticky=W, padx=5, pady=5)
48
49 myWindow.mainloop()
50
```

实验总结：

- 1.学习了 python 的基本语句，了解了 python 的使用方法和常规函数。
- 2.掌握了爬虫的基础原理和简单实用。
- 3.熟练了 GUI 图形交互界面的简单使用。