# 现代程序设计第14周作业

谢奕飞 20377077

## 代码

### WriteHead函数

```python
def WriteHead(head,dir_path):
    '''
    写入表头
    '''
    if not os.path.exists(dir_path):
        os.mkdir(dir_path)
    with open(dir_path+'/data.csv','w',encoding='utf8',newline='') as f:
        writer=csv.writer(f)
        writer.writerow(head)
```

### GetPage函数

```python
def GetPage(url,headers):
    '''
    获取页数
    '''
    response=requests.get(url=url,headers=headers)
    soup=BeautifulSoup(response.text,'lxml')
    return int(soup.select('a[class="zpgi"]')[-1].get_text())
```

### DownloadImage函数

```python
def DownloadImage(url,dir_path):
    '''
    下载图片
    '''
    if not os.path.exists(dir_path):
        os.mkdir(dir_path)
    name=url.split('/')[-1]
    img_path=dir_path+'/'+name
    try:
        urllib.request.urlretrieve(url,filename=img_path)
        urllib.request.urlcleanup()
    except:
        return 'error'
    return img_path
```

## Producer函数

```python
def Producer(q:Queue,url,headers):
    '''
    生产者
    '''
    response=requests.get(url=url,headers=headers)
    soup=BeautifulSoup(response.text,'lxml')
    soup_list=soup.select('a[class="tit f-thide s-fc0"]')
    href_list=[]
    for s in soup_list:
        href_list.append('https://music.163.com'+str(s['href']))
    q.put(href_list)
```

## Consumer函数

```python
def Consumer(q:Queue,headers,dir_path):
    '''
    消费者
    '''
    urls=q.get()
    if urls!=None:
        for url in urls:
            response=requests.get(url=url,headers=headers)
            soup=BeautifulSoup(response.text,'lxml')
            id=url.split('id=')[-1]#歌单id
            title=soup.select('.tit')[0].get_text()[1:]#歌单标题
            image_url=soup.select('img[class="j-img"]')[0]['data-src']#封面u
            image=DownloadImage(image_url,dir_path+'/images')#下载图片
            author_id=soup.select('a[class="s-fc7"]')[0]['href'].split('id='
            author=soup.select('a[class="s-fc7"]')[0].get_text()#作者名
            description=soup.select('p')[1].get_text()#简介
            count=soup.select('span[id="playlist-track-count"]')[0].get_text
            play=soup.select('strong[id="play-count"]')[0].get_text()#播放
            add=soup.select('a[class="u-btni u-btni-fav"]')[0]['data-count']
            share=soup.select('a[class="u-btni u-btni-share"]')[0]['data-cou
            comment=soup.select('span[id="cnt_comment_count"]')[0].get_text(
            #按行写入
            with open(dir_path+'/data.csv','a',encoding='utf8',newline='') a
                writer=csv.writer(f)
                writer.writerow([id,title,image,author_id,author,descriptio
```
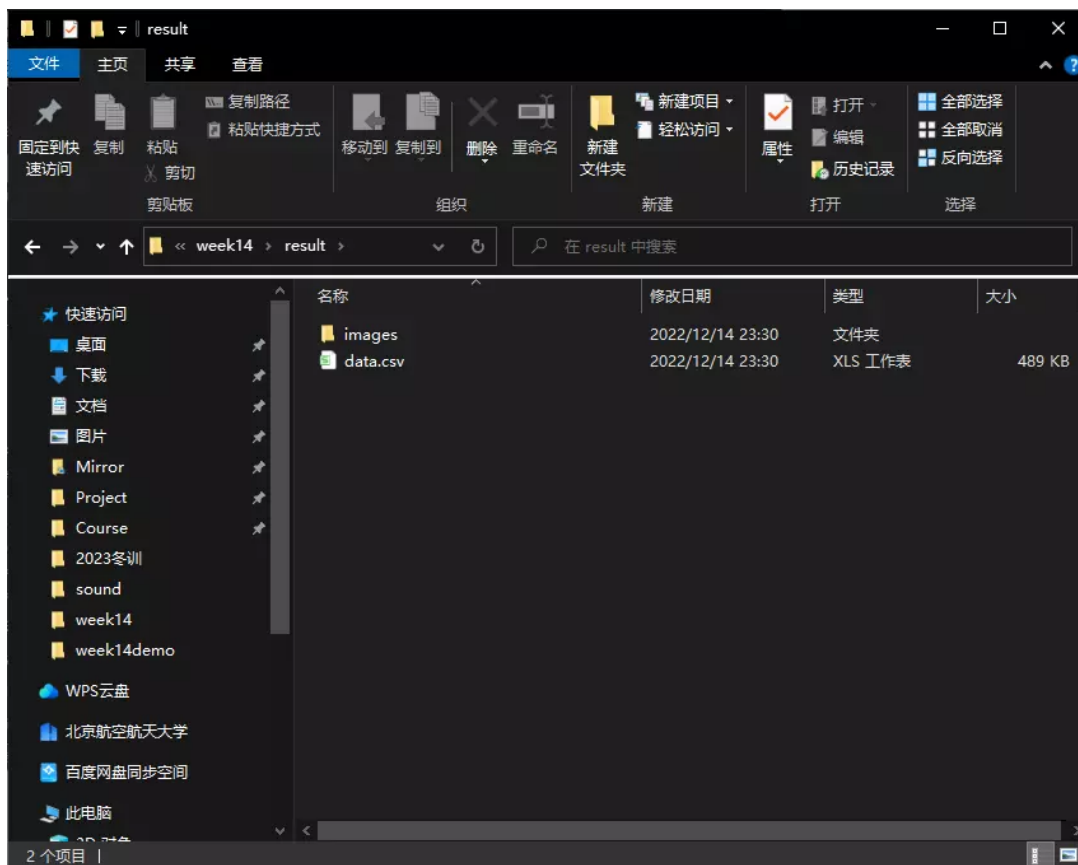
## CoroProducer函数

```python
def CoroProducer(q:Queue,url):
    '''
    协程执行Producer
    '''
```

```
 5      tasks=[]
 6      headers={'User-Agent':'Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWe
 7      N=GetPage(url,headers)#获取页数
 8      print('This website has {} pages to crawl'.format(N))
 9      #N=3
10      urls=[]
11      for i in range(N):
12          urls.append(url[:-1]+str(i*35))
13      for url in urls:
14          task=gevent.spawn(Producer,q,url,headers)
15          tasks.append(task)
16      gevent.joinall(tasks)
17      return N
```

## CoroConsumer函数

```
 1  def CoroConsumer(q:Queue,N,dir_path):
 2      '''
 3      协程执行Consumer
 4      '''
 5      headers={'User-Agent':'Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWe
 6      tasks=[]
 7      for i in range(N):
 8          task=gevent.spawn(Consumer,q,headers,dir_path)
 9          tasks.append(task)
10      gevent.joinall(tasks)
```

## main函数

```
 1  if __name__=='__main__':
 2      dir_path='D:/Project/Python/week14/result'
 3      head=['id','title','image','author_id','author','description','count','p
 4      url='https://music.163.com/discover/playlist/?order=hot&cat=%E6%B0%91%E8
 5      q=Queue()
 6      t_start=time.time()
 7      N=CoroProducer(q,url)#协程执行Producer
 8      print('Producers have been executed, all tasks cost {}s'.format(time.tim
 9      WriteHead(head,dir_path)#写入表头
10      CoroConsumer(q,N,dir_path)#协程执行Consumer
11      print('Consumers have been executed, all tasks cost {}s'.format(time.tim
12      playsound('D:\Project\Python\week14\over.mp3')
```

# 运行结果

This website has 44 pages to crawl
Producers have been executed, all tasks cost 26.43780493736267s
Consumers have been executed, all tasks cost 1953.8772449493408s
请按任意键继续. . .

分别用时26.43s和1953.87s

results目录



images目录

data.csv

# 附录–完整代码

```python
import requests,time,csv,os,urllib.request,gevent
from bs4 import BeautifulSoup
from queue import Queue
from threading import Thread
from playsound import playsound

def WriteHead(head,dir_path):
    '''
    写入表头
    '''
    if not os.path.exists(dir_path):
        os.mkdir(dir_path)
    with open(dir_path+'/data.csv','w',encoding='utf8',newline='') as f:
        writer=csv.writer(f)
        writer.writerow(head)

def GetPage(url,headers):
    '''
    获取页数
    '''
    response=requests.get(url=url,headers=headers)
    soup=BeautifulSoup(response.text,'lxml')
    return int(soup.select('a[class="zpgi"]')[-1].get_text())

def DownloadImage(url,dir_path):
    '''
    下载图片
    '''
    if not os.path.exists(dir_path):
        os.mkdir(dir_path)
    name=url.split('/')[-1]
    img_path=dir_path+'/'+name
```

```python
33          try:
34              urllib.request.urlretrieve(url,filename=img_path)
35              urllib.request.urlcleanup()
36          except:
37              return 'error'
38          return img_path
39
40  def Producer(q:Queue,url,headers):
41      '''
42      生产者
43      '''
44      response=requests.get(url=url,headers=headers)
45      soup=BeautifulSoup(response.text,'lxml')
46      soup_list=soup.select('a[class="tit f-thide s-fc0"]')
47      href_list=[]
48      for s in soup_list:
49          href_list.append('https://music.163.com'+str(s['href']))
50      q.put(href_list)
51
52  def Consumer(q:Queue,headers,dir_path):
53      '''
54      消费者
55      '''
56      urls=q.get()
57      if urls!=None:
58          for url in urls:
59              response=requests.get(url=url,headers=headers)
60              soup=BeautifulSoup(response.text,'lxml')
61              id=url.split('id=')[-1]#歌单id
62              title=soup.select('.tit')[0].get_text()[1:]#歌单标题
63              image_url=soup.select('img[class="j-img"]')[0]['data-src']#封
面url
64              image=DownloadImage(image_url,dir_path+'/images')#下载图片
65              author_id=soup.select('a[class="s-fc7"]')[0]['href'].split('id
=')[-1]#作者id
66              author=soup.select('a[class="s-fc7"]')[0].get_text()#作者名
67              description=soup.select('p')[1].get_text()#简介
68              count=soup.select('span[id="playlist-track-count"]')[0].get_te
xt()#歌曲数
69              play=soup.select('strong[id="play-count"]')[0].get_text()#播放
次数
70              add=soup.select('a[class="u-btni u-btni-fav"]')[0]['data-coun
t']#收藏数
71              share=soup.select('a[class="u-btni u-btni-share"]')[0]['data-c
ount']#分享数
72              comment=soup.select('span[id="cnt_comment_count"]')[0].get_tex
t()#评论数
73              #按行写入
```

```python
74                with open(dir_path+'/data.csv','a',encoding='utf8',newline='')
     as f:
75                    writer=csv.writer(f)
76                    writer.writerow([id,title,image,author_id,author,descripti
     on,count,play,add,share,comment])
77
78   def CoroProducer(q:Queue,url):
79       '''
80       协程执行Producer
81       '''
82       tasks=[]
83       headers={'User-Agent':'Mozilla/5.0 (Windows NT 10.0; Win64; x64) Apple
     WebKit/537.36 (KHTML, like Gecko) Chrome/107.0.0.0 Safari/537.36 Edg/107.
     0.1418.35'}
84       N=GetPage(url,headers)#获取页数
85       print('This website has {} pages to crawl'.format(N))
86       #N=3
87       urls=[]
88       for i in range(N):
89           urls.append(url[:-1]+str(i*35))
90       for url in urls:
91           task=gevent.spawn(Producer,q,url,headers)
92           tasks.append(task)
93       gevent.joinall(tasks)
94       return N
95
96   def CoroConsumer(q:Queue,N,dir_path):
97       '''
98       协程执行Consumer
99       '''
100      headers={'User-Agent':'Mozilla/5.0 (Windows NT 10.0; Win64; x64) Apple
     WebKit/537.36 (KHTML, like Gecko) Chrome/107.0.0.0 Safari/537.36 Edg/107.
     0.1418.35'}
101      tasks=[]
102      for i in range(N):
103          task=gevent.spawn(Consumer,q,headers,dir_path)
104          tasks.append(task)
105      gevent.joinall(tasks)
106
107  if __name__=='__main__':
108      dir_path='D:/Project/Python/week14/result'
109      head=['id','title','image','author_id','author','description','coun
     t','play','add','share','comment']
110      url='https://music.163.com/discover/playlist/?order=hot&cat=%E6%B0%91%
     E8%B0%A3&limit=35&offset=0'
111      q=Queue()
112      t_start=time.time()
113      N=CoroProducer(q,url)#协程执行Producer
```

```
114    print('Producers have been executed, all tasks cost {}s'.format(time.t
ime()-t_start))
115    WriteHead(head,dir_path)#写入表头
116    CoroConsumer(q,N,dir_path)#协程执行Consumer
117    print('Consumers have been executed, all tasks cost {}s'.format(time.t
ime()-t_start))
118    playsound('D:\Project\Python\week14\over.mp3')
```