

20377199 赵芮箐 第12周作业

作业内容：以爬取网易云歌单为例，练习多线程的使用。

1. 获取一个分类下的所有歌单的id。观察url可以发现其页码规律：<https://music.163.com/#/discover/playlist?order=hot&cat=%E8%AF%B4%E5%94%B1&limit=35&offset=35>。offset是本页开始的数据位置，第1页是0，第2页是35。分类参数是utf-8编码后的汉字“说唱”，使用str.encode可以处理。
2. 对每个id，获取歌单的详细信息，至少包括：歌单的封面图片（需把图片保存到本地）、歌单标题、创建者id、创建者昵称、介绍、歌曲数量、播放量、添加到播放列表次数、分享次数、评论数。可以自行实现其他信息的获取。基本信息汇总到同一张表中，以csv文件保存。<https://music.163.com/playlist?id=3037221581>
3. 要求使用生产者-消费者模式实现，要求1作为生产者，每次请求后产生新的任务交给消费者，消费者执行要求2。
4. (附加) 爬虫程序往往需要稳定运行较长的时间，因此如果你的程序突然中断或异常（比如网络或被封），如何能够快速从断点重启？
5. (附加) 爬虫程序往往需要比较好的状态输出，因此可否专门有一个线程 动态地进行输出更新，来显示当前的状态，比如程序连续运行的时长，要完成的 总页面数，其中有多少已被爬取，已收集的文件占用了多少空间，大概还需要多少时间才能完成，预计需要耗费多少硬盘空间等。

Part1: 获取一个分类下所有歌单的id

```
def producer(q, url):
    """
    生产者：获取一个分类下的所有歌单的id
    """
    headers = {
        'User-Agent': 'Mozilla/5.0 (Windows NT 6.1; WOW64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/63.0.3239.132 Safari/537.36'
    }
    response = requests.get(url=url, headers=headers)
    html = response.text
    soup = BeautifulSoup(html, 'html.parser')
    ids = soup.select('.dec a')          # 获取包含歌单详情页网址的标签
    q.put(ids)
```

Part2: 对每个id获取歌单的详细信息

```
def consumer(q):
    """
    消费者：对每个id获取歌单的详细信息
    """
    headers = {
        'User-Agent': 'Mozilla/5.0 (Windows NT 6.1; WOW64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/63.0.3239.132 Safari/537.36'
    }
    while True:
        ids = q.get()
        if ids is None:
            break
        else:
            for i in ids:
                url = 'https://music.163.com/' + i['href']          # 获取歌单详情页地址
                response = requests.get(url=url, headers=headers)
                html = response.text
                soup = BeautifulSoup(html, 'html.parser')
                # 获取歌单标题
                title = soup.select('h2')[0].get_text().replace(' ', ', ')
                # 获取封面图片
                img = soup.select('img')[0]['data-src']
```

```

res = requests.get(img)
image = Image.open(BytesIO(res.content))
image.save('week12/cover/'+ title + '.jpg')
# 获取创建者id
idd = soup.select('.s-fc7')[0]['href'].split('=')[-1]
# 获取创建者昵称
nickname = soup.select('.s-fc7')[0].get_text()
# 获取歌单介绍
if soup.select('#album-desc-more'):
    text = soup.select('#album-desc-more')[0].get_text().replace('\n', '').replace(',', ' ', ' ')
else:
    text = '无'
# 歌单内歌曲数
songs = soup.select('#playlist-track-count')[0].get_text()
# 播放量
plays = soup.select('.s-fc6')[0].get_text()
# 添加进播放列表次数
adds = soup.find('a', 'u-btni u-btni-fav')['data-count']
# 分享次数
shares = soup.find('a', 'u-btni u-btni-share')['data-count']
# 歌单评论数
comments = soup.select('#cnt_comment_count')[0].get_text().replace('评论', '0')
#print(title,idd,nickname,text,songs,plays,adds,shares,comments)
# 将详情页信息写入CSV文件中
with open('week12/music_message.csv', 'a+') as f:
    f.write(title + ',' + idd + ',' + nickname + ',' + text +
            ',' + songs + ',' + plays + ',' + adds + ',' + shares + ',' + comments + '\n')

```

Part3: 主函数实现生产者-消费者模式

```

if __name__=='__main__':
    head = ['歌单标题', 'id', '昵称', '介绍', '歌曲数量', '播放量', '添加到播放列表次数', '分享次数', '评论数']
    with open('week12/music_message.csv', 'w', encoding='utf-8') as f:
        csv_writer = csv.writer(f)          # csv格式写入文件file
        csv_writer.writerow(head)

    q = Queue()
    plist = []
    clist = []
    for i in range(0, 1300, 35):
        time.sleep(2)
        url = 'https://music.163.com/discover/playlist/?cat=后摇&order=hot&limit=35&offset=' + str(i)
        p = Thread(target=producer, args=(q, url))
        plist.append(p)
    for p in plist:
        p.start()
    for t in plist:
        p.join()
    for i in range(30):
        c = Thread(target=consumer, args=(q,))
        clist.append(c)

```

```

for c in clist:
    c.start()
for c in clist:
    q.put(None)

```

Part4: 如果你的程序突然中断或异常，如何能够快速从断点重启（附加）

修改producer代码如下：

```

try:
    response = requests.get(url=url, headers=headers)
    html = response.text
    soup = BeautifulSoup(html, 'html.parser')
    ids = soup.select('.dec a')          # 获取包含歌单详情页网址的标签
    q.put(ids)
except:
    print(f"Error, url: {url}")

```

我实现的想法比较简单，因为爬取网易云歌单的任务抓取的url都是有规律的，所以因为特殊情况停止之后可以直接从下一个url抓取的。所以在每次停止的时候我记录一下断点即可，下一次重启的时候就从该url再开始。

对于更大规模或更复杂的任务会有更好and更框架的写法，简单了解了一下，见Ref

Part5: 实现一个线程显示当前爬取状态（附加）

可以查看程序连续运行的时长，要完成的总页面数，其中有多少已被爬取，已收集的文件占用了多少空间，大概还需要多少时间才能完成，预计需要耗费多少硬盘空间等。

```

# 实时显示当前爬取状态的函数
def status():
    start = time.time()
    while True:
        time.sleep(0.5)
        now_time = time.time()-start
        with open('week12/music_message.csv', 'rb') as f:
            length = len(f.readlines()) - 1
            size = os.path.getsize('week12/music_message.csv')
            end_time = now_time / (length + 1) * 1300
            end_size = size / (length + 1) * 1300

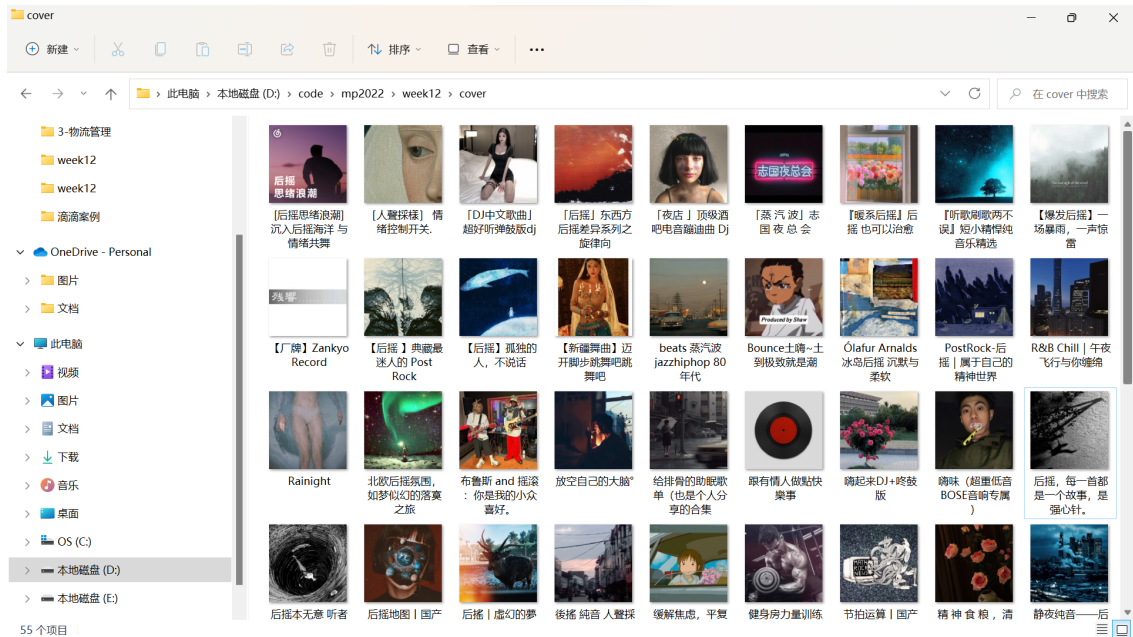
        print(f'\r运行时间: {now_time: .4f} / {end_time: .4f} \
              爬取页数: {length} / 1300 \
              文件内存: {size} / {end_size: .0f}', end='')

# 并在main里设置成守护线程
if __name__=='__main__':
    s = Thread(target=status)
    s.daemon = True
    s.start()

```

结果展示:

- 封面图片：



- 歌单信息:

