现代程序设计第十四次作业

所有作业上传至github,地址为'https://github.com/hxx12138/python_mp2021_project.git'

作业要求

以Bilibili热榜为例,练习协程与非关系数据库的使用

库的导入

```
import pymongo
from pymongo import collection
import pandas as pd
import datetime
import pickle
import os
import time
import requests as rs
import asyncio
import aiohttp
import random
from email.header import Header
from email.mime.text import MIMEText
from smtplib import SMTP_SSL
```

作业要求1

获取某一分区视频的排行榜信息,并从中解析出前 10000 条视频信息。可使用协程爬取数据。注意协程应用的方式:发送http请求后,可以跳转到数据处理的函数中,而不要跳到发送另一个http请求的函数中,因为由于b站的风控较为严格,短时间大量的请求可能会导致ip被封禁。

```
# 通过修改学长的代码实现爬虫

class BiliHotCrawler():
    def __init__(self,cate,limit=10000):
        self.cate_id = cate
        self.limit = limit
        self.page = 1
        self.pagesize = 100
        self.bulid_param()
        self.url = BASE_URL

def bulid_param(self):
```

```
self.params = {
            'main ver':'v3',
            'search type':'video',
            'view type': 'hot rank',
            'order':'click',
            'cate id':self.cate id,
            'page':self.page,
            'pagesize':self.pagesize,
            'time from':20211117,
            'time to':20211124
        }
   async def get_resp(self):
        time.sleep(random.choice(sleep choice))
        self.bulid param()
        async with aiohttp.ClientSession() as session:
            async with session.get(url=self.url,params=self.params) as response:
                self.resp = await response.read()
        self.page += 1
   def save_resp(self):
        path = f"{SAVE PATH}/{self.cate id}.csv"
        with open(path, "a") as f:
            self.resp = eval(self.resp.decode('utf-
8').replace('null','None').replace('false',"False").replace("true","True").replace("\n"
," "))
            for item in self.resp["result"]:
                title = str(item["title"]).replace('\n','o').replace(',','')
                description = str(item["description"]).replace('\n','.o').replace(',','
')
                # 分类 排名 bv号 时长 播放量 弹幕 标题 封面 评论 收藏 描述 直链
                f.write(f'{self.cate id},{item["rank offset"]},{item["bvid"]},
{item["duration"]},{item["play"]},{item["video_review"]},{title},{item["pic"]},
{item["review"]}, {item["favorites"]}, {description}, {item["arcurl"]}'.replace("\n","
").encode('utf-8', 'replace').decode('utf-8'))
                f.write("\n")
   def log(self,isend=False):
        print(f"\rCrawlering {self.cate_id} {(self.page -
1)*self.pagesize}/{self.limit}",end="")
       if isend:
           print()
   async def start(self):
        while self.pagesize * self.page <= self.limit:</pre>
            self.log()
            await self.get_resp()
            self.save resp()
        self.log(isend=True)
```

```
class BHCFactory():
   def __init__(self,cate_list):
        self.cate_list = cate_list
   def produce(self): # 相当于开了多个协程去分别爬取每个板块的热榜
        loop = asyncio.get_event_loop()
        tasks = [BiliHotCrawler(i).start() for i in self.cate list]
        loop.run_until_complete(asyncio.wait(tasks))
class BiliNoticeMail:
   def __new__(cls, *args, **kwargs):
        if not hasattr(BiliNoticeMail, '_instance'):
           BiliNoticeMail._instance = object.__new__(cls, *args, **kwargs)
        return BiliNoticeMail. instance
   def __init__(self):
        self.server = SMTP SSL(SMTP SERVER, PORT)
        self.message = None
        pass
   def build ready message(self):
        self.message = MIMEText('下载完成', 'plain', 'utf-8')
        self.message['From'] = Header(FROM_ADDR, 'utf-8')
        self.message['To'] = Header(TO ADDR, 'utf-8')
        self.message['Subject'] = Header('BiliReminder:Ready', 'utf-8')
   def notice(self):
        if self.message is None:
            raise ValueError('Message is none!')
        self.server.login(FROM ADDR, PASSWD)
        self.server.sendmail(FROM_ADDR, [TO_ADDR], self.message.as_string())
        self.message = None
        self.server.quit()
    @classmethod
   def send_ready_mail(cls):
        noticer = cls()
        noticer.build_ready_message()
        noticer.notice()
        print("发信成功,请查收。")
        return noticer
def main():
   BHCFactory(CATE_LIST).produce()
   BiliNoticeMail.send_ready_mail()
```

```
if __name__ == '__main__':
    main()
```

作业要求2

将视频信息存入MongoDB数据库,并记录此视频当时的排名和数据创建时间。

```
client =
pymongo.MongoClient("mongodb+srv://hxx:hexihexiang2000@cluster0.amu8k.mongodb.net/stude
nt?retryWrites=true&w=majority")
db = client['info']
collections_1 = db['week_1']
collections 2 = db['week 2']
db_list = client.list_database_names()
print(db_list)
data list 1 = [{'分区号':0, '排名':0, 'BV号':0, '时长':0, '播放量':0, '弹幕':0, '标题':0, '封
面':0, '评论':0, '收藏':0, '描述':0, '直链':0, '时间':0} for i in range(10001)]
df_1 = pd.read_csv('week_1/21.csv')
title_1 = list(df_1)
print(title_1)
for i in range(len(title 1)):
   print(f'the no.{i} col has completed.')
   for j in range(len(list(df 1[title 1[i]]))):
       data list 1[j][title 1[i]] = list(df 1[title 1[i]])[j]
#print(data_list_1)
for i in range(len(data_list_1)):
   print(f'insert no.{i+1}')
   date = datetime.datetime.now()
   data list 1[i]['时间'] = date
   doc = collections 1.insert one(data list 1[i])
x = collections 1.find one()
print(x)
data_list_2 = [{'分区号':0, '排名':0, 'BV号':0, '时长':0, '播放量':0, '弹幕':0, '标题':0, '封
面':0, '评论':0, '收藏':0, '描述':0, '直链':0, '时间':0} for i in range(10001)]
df_2 = pd.read_csv('week_2/21.csv')
title 2 = list(df 2)
print(title_2)
for i in range(len(title_2)):
   print(f'the no.{i} col has completed.')
   for j in range(len(list(df_2[title_2[i]]))):
```

```
data_list_2[j][title_1[i]] = list(df_2[title_2[i]])[j]
#print(data_list_2)

for i in range(len(data_list_2)):
    print(f'insert no.{i+1}')
    date = datetime.datetime.now()
    data_list_2[i]['时间'] = date
    doc = collections_2.insert_one(data_list_1[i])
x = collections_2.find_one()
print(x)
```

作业要求3

爬取同一分区一周之后的热榜,并从中解析出前 10000 条视频信息。对比两次结果并更新数据库,要求第一次结果需从数据库中取出。根据以下规则更新数据库:对于仅在第二次排行榜中的视频,将信息存入数据库;对于仅在第一次排行榜中的视频,将信息从数据库中删除;对于两次排行榜都存在的视频,更新数据库,保留创建时间,增加更新时间字段,并更新排名。

```
week 1 list = collections 1.find()
'''with open('week 1/week 1 origin.pkl','wb') as f:
   pickle.dump(list(week 1 list),f)'''
i = 0
week_2_list = collections_2.find()
for info 2 in week 2 list:
   print(f'the no.{i+1} has completed.')
   i += 1
   bv = info 2['BV号']
   query = {'BV号':bv}
   if (collections_1.find(query)):
        update_dict = {'$set':collections_1.find(query)[0]}
        collections 1.update one(query,update dict)
    else:
        collections_1.insert_one(info_2)
week_1_list = collections_1.find()
'''with open('week_1/week_1_update.pkl','wb') as f:
   pickle.dump(list(week_1_list),f)'''
```

代码测试

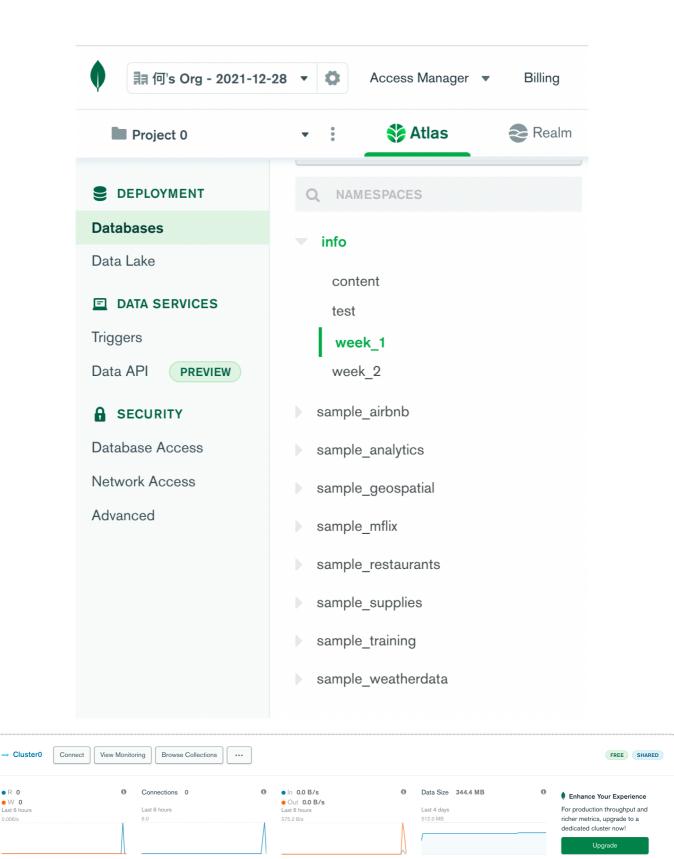
```
import pymongo
from pymongo import collection

client =
pymongo.MongoClient("mongodb+srv://hxx:hexihexiang2000@cluster0.amu8k.mongodb.net/stude
nt?retryWrites=true&w=majority")

db = client['info']
collections_1 = db['week_1']
#collections_2 = db['week_2']

doc_list_1 = collections_1.find().sort('排名')
for doc in doc_list_1:
    if doc['排名'] != 0:
        print(doc)
```

运行结果



BACKUPS

Inactive

Replica Set - 3 nodes

LINKED REALM APP

None Linked

ATLAS SEARCH

Create Index

• R 0

VERSION

4.4.10

REGION

GCP / Taiwan (asia-east1)

CLUSTER TIER

M0 Sandbox (General)

