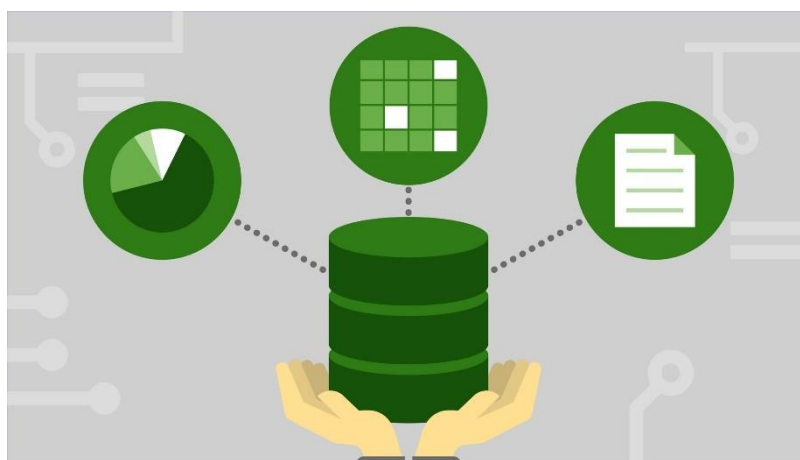


به نام خدا



دانشگاه تهران
پردیس دانشکده‌های فنی
دانشکده برق و کامپیوتر



آزمایشگاه پایگاه داده

دستور کار شماره ۶

نام و نام خانوادگی

معین شیردل ۸۱۰۱۹۷۵۳۵

آبان ماه ۱۴۰۰

گام اول- import کردن داده در دیتابیس Mongo:

ابتدا به صورت دستی و با دستور `db.tweets.save(tweet_i)` تا توییت به صورت دستی ذخیره شد و سپس، به خاطر نیاز به داشتن دقیقاً ۵۰۰ توییت، این ۱۰ توییت با ۵۰۰ توییت موجود حاصل از کد پایتونی زیر جایگزین شد.

```
tweet_extraction.py
1 from pymongo import MongoClient
2 import requests, time
3
4 tweets_url = "https://www.sahamyab.com/guest/twiter/list?v=1.0"
5 client = MongoClient('localhost', 27017)
6 db = client.Lab6
7 tweets = db.tweets
8
9 while (tweets.count_documents({}) < 500):
10     response = requests.request('GET', tweets_url, headers={ 'User-Agent': 'Chrome/61'})
11     if response.status_code == requests.codes.ok:
12         crawled_tweets = response.json()['items']
13         for tw in crawled_tweets:
14             print(tw)
15             try:
16                 tweets.replace_one({"id": tw["id"]}, tw, upsert=True)
17                 print(tw)
18             except Exception as e:
19                 print("Error occured: ", str(e))
20             time.sleep(60)
21
22 print("Count: " % tweets.count_documents({}))
```

در نهایت، شکل کلی کالکشن توییت ها به شکل زیر درآمد که قرار داشتن ۵۰۰ توییت در آن مشخص است.

The screenshot displays the MongoDB Compass interface for the 'Lab6.tweets' collection. The left sidebar shows the database structure with 'Lab6' and its collections. The main area shows a list of documents. The first document is expanded, showing fields like `_id`, `id`, `sendTime`, `sendTimePersian`, `parentSendTime`, `parentSendTimePersian`, `parentId`, `parentSenderId`, `parentSenderName`, `parentSenderProfileImage`, `parentContent`, `senderName`, `senderUsername`, `senderProfileImage`, `content`, `type`, `scoreDestDate`, and `finalPollDatePersian`. The status bar at the bottom indicates 500 documents, 497.0KB total size, and 994B avg size.

گام دوم - پیش پردازش داده:

```

from pymongo import MongoClient
import requests, time, re # re for regex

mongodb_cli = MongoClient('localhost', 27017)

db = mongodb_cli.lab6
tweets = db.tweets
start = time.time()

update_count = 0
for tw in tweets.find({}):
    tweet_content = tw['content']
    hashtags = list(re.findall(r'#(\w+)', tweet_content))
    tweets.update_one({"id": tw['_id']}, {"$set": {"hashtags": hashtags}})
    update_count += 1

end = time.time()
print("{} rows updated in {} seconds".format(update_count, end - start))

```

500 rows updated in 0.34645795822143555 seconds

در این مرحله، به کمک تکه کد بالا برای هر توییت از فیلد content آن هشتگ ها استخراج شدند و همانطور که خواسته شده بود، در فیلدی به نام hashtags ذخیره شدند. عمل اصلی در خط مشخص شده در این تکه کد انجام میگیرد که به کمک regex تمامی هشتگ های موجود را در یک لیست می ریزد. البته زمان اجرای اولیه کمی بیشتر از این مقدار بود و ۰/۷۵ ثانیه به طول انجامید که از آن تصویر تهیه نشد و این اجرای دوم، به علت تکراری بودن سریع تر انجام شد. در قسمت زیر، یک سند نمونه را می بینیم که هر ۳ هشتگ موجود در متن آن در لیست قرار دارند.

```

content: "شاخص بورس #شستا #خودرو"
...اگر کسی یک دستاورد به جز دکتر کردن یک شبه ری
lastLikeNickName: "ili"
likeCount: "45"
commentCount: "1"
retwitCount: "4"
quoteCount: "2"
type: "retwit"
scoredPostDate: "1639562834417"
retwitId: "405695598"
finalPullDatePersian: ""
▼ hashtags: Array
  0: "شاخص_بورس"
  1: "شستا"
  2: "خودرو"

```

گام سوم- دستورات اصلی:

۱) پیدا کردن نام کاربرانی که parentId دارند و mediaContentType شان image/jpeg است:

```
_id: ObjectId("61baf5f05102f18e51b38984")
parentId: "405754996"
senderName: "trader"
mediaContentType: "image/jpeg"
```

```
_id: ObjectId("61baf5f05102f18e51b3898a")
parentId: "405720345"
senderName: "AAAA"
mediaContentType: "image/jpeg"
```

```
_id: ObjectId("61bafcf85102f18e51b3898f")
parentId: "405724197"
senderName: "سینا مالکی"
mediaContentType: "image/jpeg"
```

```
_id: ObjectId("61bb191b5102f18e51b389b8")
parentId: "404995520"
senderName: "shahramdehqani"
mediaContentType: "image/jpeg"
```

Query Performance Summary

Documents Returned: 22

Actual Query Execution Time (ms): 1

FILTER {parentId: {\$exists: true}, mediaContentType: "image/jpeg"}

PROJECT {senderName: 1, mediaContentType: 1, parentId: 1}

SORT { field: -1 } or [['field', -1]]

COLLATION { locale: 'simple' }

به کمک فیلتر هایی که روی فیلدهای parentId و mediaContentType شد و projection روی سه فیلد مد نظر، نتیجه به شکلی که در تصویر سمت چپ مشخص است به نمایش در آمد. در مجموع ۲۲ رکورد این شرایط را دارا بودند.

۲) مقدار senderUsername و type از توییت های ارسالی در یک بازه ۱۵ دقیقه ای دلخواه:

FILTER {sendTime: {\$gt: "2021-12-14T21:00:00Z", \$lte: "2021-12-14T21:15:00Z"}} **OPTIONS** **FIND**

PROJECT {senderUsername:1, type:1, sendTime:1}

SORT { field: -1 } or [['field', -1]] **MAX TIME MS** 60000

COLLATION { locale: 'simple' } **SKIP** 0 **LIMIT** 0

VIEW **3 of 3**

Displaying documents 1 - 3 of 3

```
_id: ObjectId("61b911165372d8c9291017f0")
sendTime: "2021-12-14T21:15:00Z"
senderUsername: "bitcoin_miner"
```

```
_id: ObjectId("61b911165372d8c9291017f1")
sendTime: "2021-12-14T21:14:33Z"
senderUsername: "leo.king"
```

```
_id: ObjectId("61b911165372d8c9291017f2")
sendTime: "2021-12-14T21:14:04Z"
senderUsername: "bourse_iran_"
type: "twit"
```

Query Performance Summary

Documents Returned: 3

Actual Query Execution Time (ms): 4

برای این سوال، بازه زمانی ساعت ۹:۰۰ تا ۹:۱۵ تاریخ ۱۴ دسامبر استفاده شد که طبق تصویر، ۱۱ توییت در این بازه زمانی منتشر شده بودند. مقادیر sendTime آن ها باید در این بازه باشد چون sendTime حاصل تبدیل DateTime به استرینگ است و برای مقایسه مناسب تر است.

۳) کاربرانی که از ساعت ۹ تا ۱۰ یک روز دلخواه بیش از یک توییت کرده اند:

```

1 FILTER {sendTime: {$gte: "2021-12-16T09:00:00Z", $lt: "2021-12-16T10:00:00Z"}}
2 PROJECT [{senderUsername: 1, senderName: 1, sendTime: 1}]
3 SORT [{senderUsername: 1}]
4 COLLATION { locale: 'simple' }

```

Query Performance Summary

Documents Returned: 20

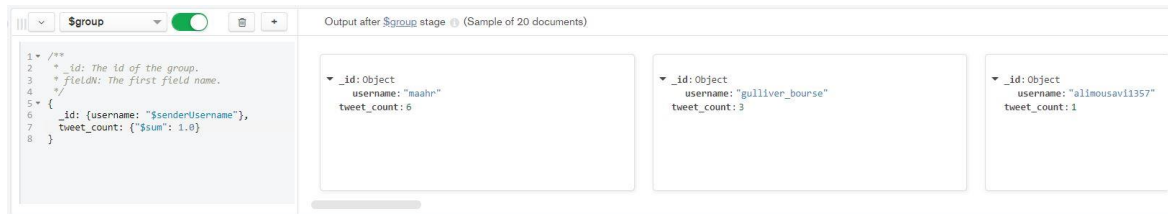
Actual Query Execution Time (ms): 3

برای این سوال، از آنجایی که به نظر باید بدون استفاده از aggregation و گروه‌بندی روی username کاربران اسامی را پیدا می‌کردیم، ساعت ۹ تا ۱۰ صبح روز ۱۶ دسامبر انتخاب شد که ۲۰ توییت در این بازه موجود بود. سپس برای راحت تر متوجه شدن تعداد توییت هر کاربر، یک sorting روی senderUsername انجام گرفت تا username های مشابه در کنار هم قرار گیرند و تعداد توییت های هر کس راحت تر به چشم آید. به طور مثال، در این ۲۰ توییت دو کاربر بودند که بیش از یک توییت داشتند و می‌توانند جایزه دریافت کنند:

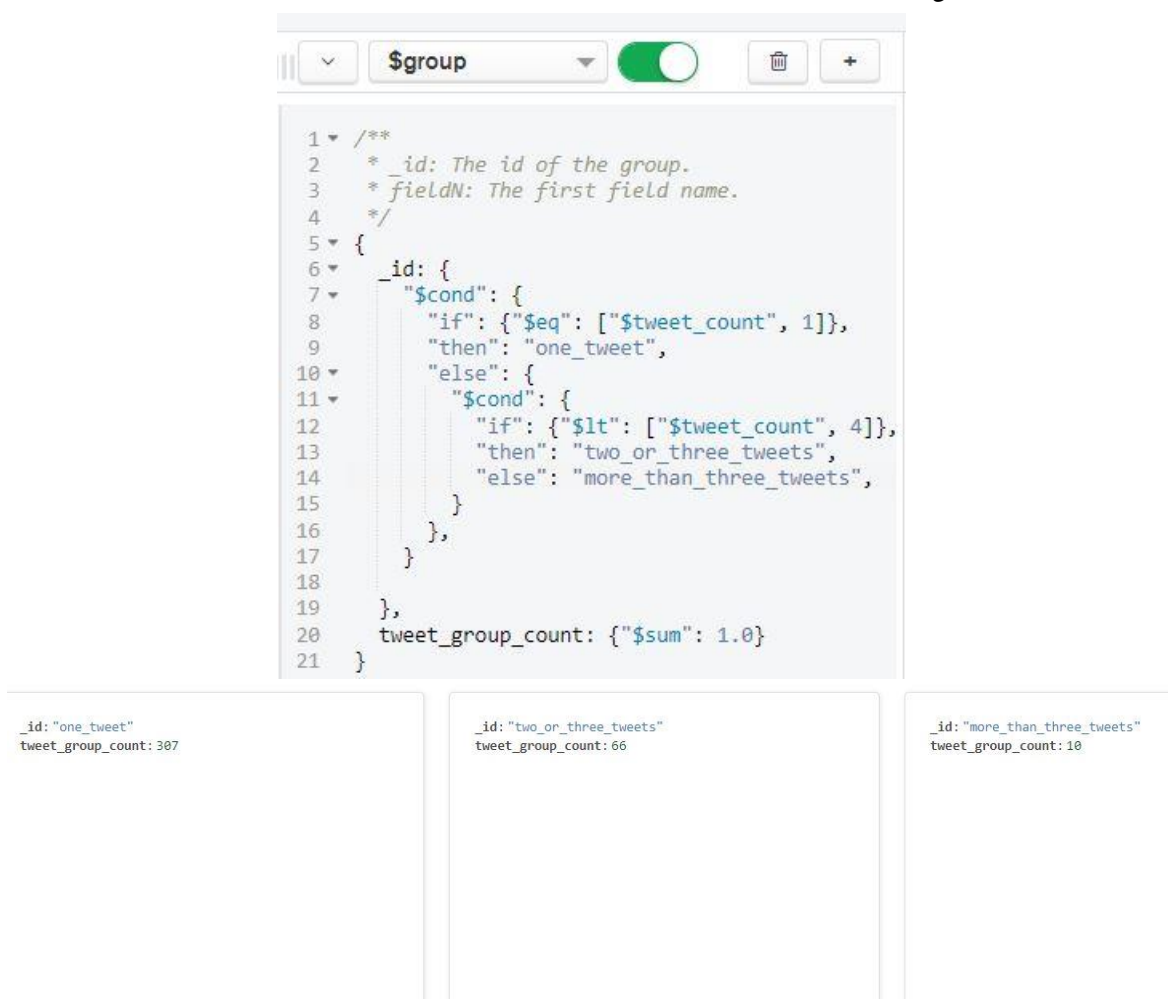
<pre> _id: ObjectId("61bb0b0a5102f18e51b389a5") sendTime: "2021-12-16T09:39:46Z" senderName: "فداتجلیل" senderUsername: "indexcode" </pre>	<pre> _id: ObjectId("61bb0b0a5102f18e51b389a2") sendTime: "2021-12-16T09:41:43Z" senderName: "علی" senderUsername: "alidsh88" </pre>
<pre> _id: ObjectId("61bb0b0a5102f18e51b389a7") sendTime: "2021-12-16T09:38:53Z" senderName: "فداتجلیل" senderUsername: "indexcode" </pre>	<pre> _id: ObjectId("61bb0b0a5102f18e51b389a0") sendTime: "2021-12-16T09:43:42Z" senderName: "علی" senderUsername: "alidsh88" </pre>

گام چهارم- دستورات تجمعی و آماری (Aggregate Functions):

۱) یافتن تعداد کاربران بر اساس گروه‌های مربوط به تعداد توییت (یک توییت، ۲ یا ۳ توییت و بیش از ۳ توییت):
به کمک محیط گرافیکی و با aggregation pipeline با stage این مورد انجام شد.



در اولین استیج، تعداد توییت‌های هر کاربر در فیلد tweet_count ذخیره شد. نمونه‌ای از نتایج پس از این مرحله نیز در سمت راست تصویر قابل مشاهده است.



تکه کد بالا نیز مربوط به استیج دوم از پایپ لاین است و خروجی آن در تصویر دوم مشخص است. کد پایتونی این دستورها نیز در ادامه آمده است.

```

aggregation_res = db.tweets.aggregate([
    {
        '$group': {
            '_id': {
                'username': '$senderUsername'
            },
            'tweet_count': {
                '$sum': 1
            }
        }
    }, {
        '$group': {
            '_id': {
                '$cond': {
                    'if': {
                        '$eq': [
                            '$tweet_count', 1
                        ]
                    },
                    'then': 'one_tweet',
                    'else': {
                        '$cond': {
                            'if': {
                                '$lt': [
                                    '$tweet_count', 4
                                ]
                            },
                            'then': 'two_or_three_tweets',
                            'else': 'more_than_three_tweets'
                        }
                    }
                }
            },
            'tweet_group_count': {
                '$sum': 1
            }
        }
    }, {
        '$sort': {
            'tweet_group_count': -1
        }
    }
])
end = time.time()
for res in aggregation_res:
    print(res)

print("Completed in {} seconds".format(end - start))

{'_id': 'more_than_three_tweets', 'tweet_group_count': 10}
{'_id': 'one_tweet', 'tweet_group_count': 307}
{'_id': 'two_or_three_tweets', 'tweet_group_count': 66}
Completed in 0.017939329147338867 seconds

```

در تصویر بالا، نتایج و تعداد کاربران هر دسته مشخص است و طول زمان اجرای کد نیز آورده شده است.

۲) یافتن پرتکرار ترین هشتکها:

```

aggregation_res = tweets.aggregate([
    {"$unwind": "$hashtags"},
    {"$group":
        {
            '_id': {'hashtag': '$hashtags'},
            'hashtag_count': {'$sum': 1}
        }
    },
    {"$sort": {"hashtag_count": -1}}
])

end = time.time()
print("Completed in {} seconds. \n\n".format(end - start))
for res in aggregation_res:
    print(res)

```

Completed in 0.007529735565185547 seconds.

```

{'_id': {'hashtag': 'شاخص بورس'}, 'hashtag_count': 44}
{'_id': {'hashtag': 'خودرو'}, 'hashtag_count': 24}
{'_id': {'hashtag': 'برکت'}, 'hashtag_count': 19}
{'_id': {'hashtag': 'شیلی'}, 'hashtag_count': 17}
{'_id': {'hashtag': 'شستا'}, 'hashtag_count': 15}
{'_id': {'hashtag': 'کسرا'}, 'hashtag_count': 13}
{'_id': {'hashtag': 'پالایش'}, 'hashtag_count': 10}
{'_id': {'hashtag': 'نیکو'}, 'hashtag_count': 8}
{'_id': {'hashtag': 'کما'}, 'hashtag_count': 7}

```

ابتدا لیست هشتک های هر توییت را unwind می کنیم. این عمل معادل با explode کردن یک لیست است که به ازای هر مقدار از لیست هشتکها، یک سند جدید ایجاد می کند و به عبارتی لیست را با المان هایش در سندهای مختلف جایگزین می کند. سپس به روش گروه بندی کردن روی نام هشتکها، تعداد تکرار هریک را می شماریم که نتایج بالا حاصل می شود. پرتکرار ترین، مطابق انتظار هشتک شاخص_بورس است. در تصویر خروجی کامل قرار داده نشده است و تعداد از پرتکرار ترین ها را می بینیم.

۳) حذف فیلد type از توییت های دارای parent-id:

```

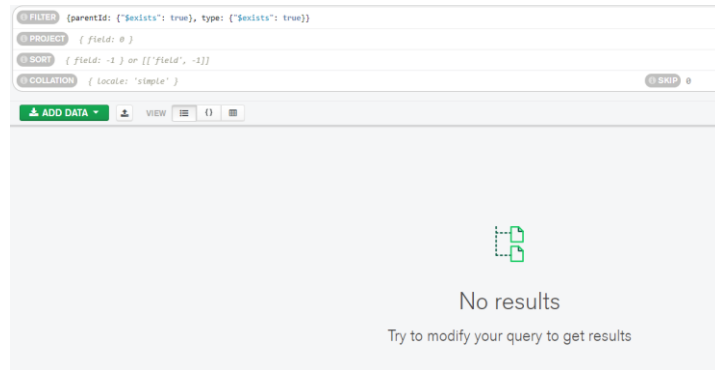
res = tweets.update_many(
    {"parentId": {"$exists": "true"}},
    {"$unset": {"type": ""}}
)

end = time.time()
print("Completed in {} seconds. \n\n".format(end - start))

```

Completed in 0.007972002029418945 seconds.

به کمک دستور unset برای رکورد های مورد نظر فیلد type را حذف می کنیم. سپس برای اطمینان از حذف شدن آن، در تصویر پایین می بینیم که آیا رکوردی هست که هم parentId داشته باشد و هم type داشته باشد یا نه. که همانگونه که در تصویر پایین مشخص است، این کوئری نتیجه ای در بر نداشته و کار ما به درستی انجام شده است.



۴) پرتکرارترین و کم‌تکرارترین هشتگ‌ها:

```
aggregation_res = tweets.aggregate([
    { "$match": {
        "sendTime": {
            "$gte": start_datetime,
            "$lte": end_datetime
        }
    },
    { "$unwind": "$hashtags" },
    { "$group": {
        "_id": { 'hashtag': '$hashtags' },
        'hashtag_count': { '$sum': 1 }
    } },
    { "$sort": { "hashtag_count": -1 } }
])

end = time.time()
print("Completed in {} seconds. \n\n".format(end - start))
aggregation_res = list(aggregation_res)

max_hashtag_count = aggregation_res[0]["hashtag_count"]
min_hashtag_count = aggregation_res[len(aggregation_res) - 1]["hashtag_count"]

print("Max hashtag count: ")
i = 0
while aggregation_res[i]["hashtag_count"] == max_hashtag_count:
    print(aggregation_res[i]["_id"], max_hashtag_count)
    i += 1

print("\n\nMin hashtag count: ")
i = len(aggregation_res) - 1
while aggregation_res[i]["hashtag_count"] == min_hashtag_count:
    print(aggregation_res[i]["_id"], min_hashtag_count)
    i -= 1
```

Completed in 0.007975339889526367 seconds.	{ 'hashtag': 'کمشه' } 1	{ 'hashtag': 'سننیر' } 1
	{ 'hashtag': 'اسموی' } 1	{ 'hashtag': 'زینبا' } 1
	{ 'hashtag': 'کفایرس' } 1	{ 'hashtag': 'خکوه' } 1
	{ 'hashtag': 'سمنها' } 1	{ 'hashtag': 'شیاکسا' } 1
	{ 'hashtag': 'سگا' } 1	{ 'hashtag': 'نوب' } 1
	{ 'hashtag': 'یکام' } 1	{ 'hashtag': 'ویانک' } 1
	{ 'hashtag': 'فجر' } 1	{ 'hashtag': 'ولسایا' } 1
	{ 'hashtag': 'خپهن' } 1	{ 'hashtag': 'وایند' } 1
	{ 'hashtag': 'خزایما' } 1	{ 'hashtag': 'امید' } 1
	{ 'hashtag': 'اسید' } 1	{ 'hashtag': 'پورس' } 1
	{ 'hashtag': 'جریل' } 1	{ 'hashtag': 'ستران' } 1
	{ 'hashtag': 'حسینا' } 1	{ 'hashtag': 'غینو' } 1
	{ 'hashtag': 'غکورش' } 1	{ 'hashtag': 'کیما' } 1
	{ 'hashtag': 'فاسم' } 1	{ 'hashtag': 'سبجنو' } 1
	{ 'hashtag': 'سکارون' } 1	{ 'hashtag': 'بالاس' } 1
	{ 'hashtag': 'تیرسا' } 1	{ 'hashtag': 'فخور' } 1
	{ 'hashtag': 'تاخت' } 1	{ 'hashtag': 'نکران' } 1
	{ 'hashtag': 'رافزا' } 1	{ 'hashtag': 'کازرو' } 1
	{ 'hashtag': 'جهیم' } 1	{ 'hashtag': 'وسید' } 1
	{ 'hashtag': 'ومستوی' } 1	{ 'hashtag': 'وغنیر' } 1
	{ 'hashtag': 'خفاری' } 1	
	{ 'hashtag': 'توریل' } 1	
	{ 'hashtag': 'دایران' } 1	
	{ 'hashtag': 'کالا' } 1	
	{ 'hashtag': 'ختر' } 1	
	{ 'hashtag': 'تیندر' } 1	
	{ 'hashtag': 'والدر' } 1	
	{ 'hashtag': 'یکینوج' } 1	

مشابه کد استفاده شده در قسمت دوم از گام چهارم، تعداد تکرار هر هشتگ را حساب می‌کنیم. سپس مقدار کمترین و بیشترین تکرار را با تبدیل حاصل به لیست استخراج می‌کنیم. در مرحله آخر برای چاپ، از ابتدای لیست شروع می‌کنیم و ماکسیمم‌ها را چاپ می‌کنیم و سپس از انتهای لیست مینیمم‌ها را.

(۵) ۱۰ هشتگ پرتکرار یک روز خاص:

```
start_datetime = "2021-12-16T00:00:00Z"
end_datetime = "2021-12-16T23:59:59Z"

aggregation_res = tweets.aggregate([
    {"$match": {
        "sendTime": {
            "$gte": start_datetime,
            "$lte": end_datetime
        }
    }},
    {"$unwind": "$hashtags"},
    {"$group": {
        '_id': {'hashtag': '$hashtags'},
        'hashtag_count': {'$sum': 1}
    }},
    {"$sort": {'hashtag_count': -1}}
])

end = time.time()
print("Completed in {} seconds. \n\n".format(end - start))
printed_count = 0
for res in aggregation_res:
    if printed_count == 10:
        break
    print(res)
    printed_count += 1
```

در کد بالا، شروع و پایان روز مدنظر (مثلاً ۱۶ دسامبر) در دو متغیر ثابت داده شده‌اند. سپس sendTime فیلتر شده است و بعد به همان روش استفاده شده در سوال دوم گام چهارم، هشتگ‌های پرتکرار شمرده شده‌اند. نتیجه در تصویر زیر قابل مشاهده است.

Completed in 0.006978273391723633 seconds.

```
{'_id': {'hashtag': 'شاخص بورس'}, 'hashtag_count': 25}
{'_id': {'hashtag': 'برکت'}, 'hashtag_count': 14}
{'_id': {'hashtag': 'خودرو'}, 'hashtag_count': 9}
{'_id': {'hashtag': 'شپلی'}, 'hashtag_count': 8}
{'_id': {'hashtag': 'پالایش'}, 'hashtag_count': 8}
{'_id': {'hashtag': 'شستا'}, 'hashtag_count': 5}
{'_id': {'hashtag': 'تشاهد'}, 'hashtag_count': 5}
{'_id': {'hashtag': 'شلرد'}, 'hashtag_count': 4}
{'_id': {'hashtag': 'کنما'}, 'hashtag_count': 4}
{'_id': {'hashtag': 'نیکو'}, 'hashtag_count': 4}
```

۶) فعال‌ترین کاربر در یک روز خاص:

```

start_datetime = "2021-12-16T00:00:00Z"
end_datetime = "2021-12-16T23:59:59Z"

aggregation_res = tweets.aggregate([
    {"$match": {
        "sendTime": {
            "$gte": start_datetime,
            "$lte": end_datetime
        }
    }},
    {"$group": {
        "_id": {'senderUsername': '$senderUsername'},
        'tweet_count': {'$sum': 1}
    }},
    {"$sort": {"tweet_count": -1}}
])

end = time.time()
print("Completed in {} seconds. \n\n".format(end - start))

aggregation_res = list(aggregation_res)
max_tweet_count = aggregation_res[0]["tweet_count"]

print("Most active user: ")
i = 0
while aggregation_res[i]["tweet_count"] == max_tweet_count:
    print(aggregation_res[i]["_id"], "\tTweets count:", max_tweet_count)
    i += 1

```

در کد بالا نیز، شروع و پایان روز مدنظر (مثلاً ۱۶ دسامبر) در دو متغیر ثابت داده شده‌اند. سپس sendTime فیلتر شده است و بعد گروه‌بندی روی senderUsername انجام گرفته و تعداد توییت‌های هر کاربر شمارش شده‌است. در پایان نیز کاربر(ها) با بیشترین تعداد توییت مشخص شده‌است.

Completed in 0.008061408996582031 seconds.

Most active user:
 {'senderUsername': 'alidsh88'} Tweets count: 8
