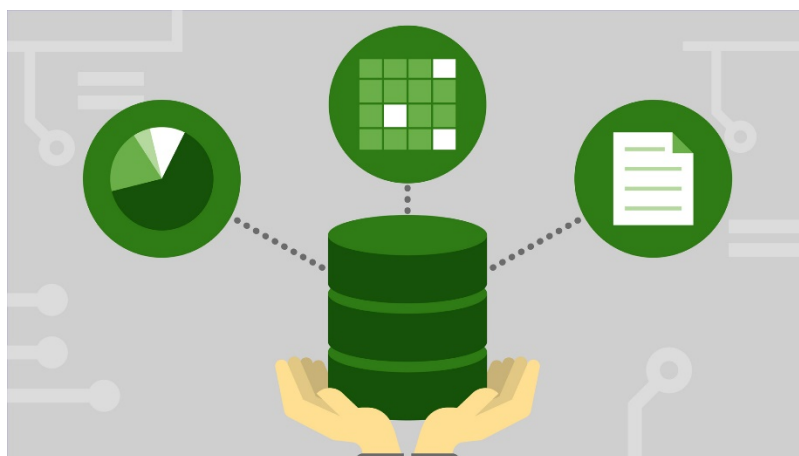


به نام خدا



دانشگاه تهران
پردیس دانشکده‌های فنی
دانشکده برق و کامپیوتر



آزمایشگاه پایگاه داده

دستورکار شماره ۶

MongoDB

مهلت تحویل: ۱۴۰۱/۱۰/۳۰

مجتبی بنائی

دستور کار شماره ۶

هدف اصلی از این تمرین، آشنایی با مانگودی بی به عنوان یکی از رایجترین دیتابیس‌های غیررابطه‌ای (NoSQL) دنیاست.

در این تمرین، برای ذخیره توییت‌های سایت سهامیاب از مانگو استفاده میکنیم و بعد از ذخیره اطلاعات، با انجام چند پرس و جوی ساده، نحوه کار با این دیتابیس محبوب را فراخواهیم گرفت.

پیش‌نیاز و شروع به کار با مانگودی بی

در این قسمت، به آشنایی با مانگودی بی به کمک فیلم‌های آموزشی سایت رسمی مانگودی بی به آدرس

<https://university.mongodb.com/courses/M001/about>

خواهیم پرداخت.

بنابراین، با مرور فیلم‌های آموزشی مانگودی بی و انجام بخشی از تمرینات انجام شده آن، از دستوراتی که نوشته‌اید و خروجی‌هایی که تولید کرده‌اید (انتخاب آنها بر عهده خودتان خواهد بود)، اسکرین شات گرفته، برای هر یک، توضیحی کوتاه بدهید. (شش کوثری و خروجی آنها برای گزارش کار کافی است.)

دقت کنید که مجموع طول دوره آموزشی فوق، هشت ساعت است اما در همان حدی از فیلم‌ها را مشاهده و به صورت عملی کار کنید که آماده انجام ادامه‌ی دستور کار شوید. (هر چند توصیه می‌کنم شش فصل این دوره را حتما مشاهده و به صورت عملی انجام دهید)

نصب و راه‌اندازی مانگودی بی

در ادامه برای یادگیری نحوه‌ی نصب و راه‌اندازی دیتابیس مانگودی بی، مقاله آشنایی با مانگو دی بی سایت مهندسی داده^۱ را مطالعه نموده و بخش پیاده سازی دیتابیس **retrogames** آنرا طبق دستوراتی که داده شده است در خط فرمان (پاورشل یا خط فرمان لینوکس) انجام دهید. سپس همین دستورات را در محیط **Robo3T** و یا **MongoDB Compass** اجرا کنید.

این بخش از کار، به عنوان پیش‌نیاز و دست‌گرمی محسوب می‌شود و نیاز به آوردن آن در گزارش نهایی دستورکار نخواهد بود.

دریافت اطلاعات

با استفاده از <https://www.sahamyab.com/guest/twitter/list?v=0.1> ده توییت آخر سایت سهامیاب با تمامی مشخصات را با فرمت جی‌سان دریافت میکنیم (با پستمن با روش GET می‌توانید خروجی را تست کنید). توییت‌ها در فیلد **items** پاسخ، قابل مشاهده هستند. برای این تمرین، به کمک API فوق به جمع‌آوری و پردازش توییت‌های فارسی خواهیم پرداخت.

برای دریافت اطلاعات می‌توانید از کد زیر استفاده کنید:

^۱ www.bigdata.ir/?p=214

```

import requests, time
url = "https://www.sahamyab.com/guest/twiter/list?v=0.1"
delay = 60
while True:
    response = requests.request('GET', url, headers={'User-Agent': 'Chrome/61'})
    if response.status_code == requests.codes.ok:
        items = response.json()['items']
        for item in items:
            print(item)
            time.sleep(60)

```

نصب مانگو و ساخت کالکشن توییت ها

مانگو دی بی را نصب کرده و کالکشن **tweets** را در دیتابیس **sahamyab** (این دیتابیس هم باید ایجاد شود) بسازید. میتوانید از خط فرمان مانگودی بی یا ابزارهای گرافیکی رایج مانند **MongoDB Compass**¹ برای این منظور استفاده کنید.

کتاب کوچک **The Little MongoDB**² می تواند راهنمای سریع شما برای کار با مانگو در این تمرین باشد.

گام اول تمرین

در این گام، با فراخوانی آدرس <https://www.sahamyab.com/guest/twiter/list?v=0.1> ده توییت آخر را دریافت کرده و به صورت دستی در مانگو ذخیره کنید (از خروجی پستمن هم می توانید در این مرحله استفاده کنید و نیاز به کدنویسی نخواهد بود) و بررسی کنید چه فیلدهایی توسط خود مانگو به صورت خودکار به داده ها افزوده میشود. (هر توییت را به عنوان یک داکيومنت ذخیره کنید یعنی با فراخوانی کد فوق، ده توییت را ذخیره خواهیم کرد.)

سپس با استفاده از کتابخانه **pymongo**³ کد دریافت اطلاعات فوق را به گونه ای تغییر دهید که هر یک دقیقه یکبار، توییت های جدید را دریافت کرده و همزمان با دریافت توییت ها، آنها را در مانگو هم ذخیره کند. (دقت کنید که هر توییت باید جداگانه ذخیره شود و توییت های تکراری بر اساس فیلد **id** هم باید حذف شوند که البته می توانید **upsert** کنید)

کد نوشته شده را تا زمانی اجرا کنید که حداقل ۵۰۰ توییت منحصر بفرد در مانگو ذخیره شده باشند. با دستور **count**، مطمئن شوید که ۵۰۰ توییت ذخیره شده باشد.

خروجی گام اول

نحوه ورود دستی داده ها در مانگو و فیلدهای اضافه شده، کدهای نوشته شده برای درج اطلاعات و نحوه اطمینان از درج ۵۰۰ توییت در گزارش آورده شود.

¹ <https://www.mongodb.com/products/compass>

² <https://openmymind.net/mongodb.pdf>

³ <https://pymongo.readthedocs.io/en/stable/tutorial.html>

گام دوم - پیش پردازش داده

در این گام با استفاده از **Regex** هشتگ های استفاده شده کاربر در فیلد **content** را پیدا کرده و سپس با استفاده از دستور **update** در فیلدی به نام **hashtags** به صورت **Array** ذخیره کنید.

خروجی گام دوم

دستور نوشته شده، خروجی و زمان اجرا

گام سوم - دستورات اصلی

1. نام کاربرانی که **mediaContentType** توییت آنها **image/png** هستند و **parentId** آنها مقدار دارد را بیابید.
2. **senderUsername** و **content** توییت آنهایی که در یک بازه ۱۵ دقیقه ای دلخواه (از بازه توییت های دریافتی (توییت فرستاده اند را بیابید).
3. قصد داریم **senderName** و **senderProfileImage** کسانی که در یک روز خاص و یک بازه ی یک ساعته که داده های آن در کالکشن شما موجود باشد و بیش از یک توییت کرده اند را پیدا کنیم، این کاربران را بیابید.

خروجی گام سوم

دستور نوشته شده، خروجی و زمان اجرا

گام چهارم - دستورات تجمعی و آماری (Aggregate Functions)

1. می خواهیم کاربران را بر اساس فعالیتشان دسته بندی کنیم. کاربران را به سه دسته به صورت زیر تقسیم کنید:
- کاربرانی با یک توییت، کاربرانی با دو تا سه توییت، کاربرانی با بیش از سه توییت
- دستوری بنویسید که تعداد هر گروه را برگرداند.
2. تعداد توییت های هر هشتگ را بشمارید و به صورت نزولی رتبه بندی کنید.
3. برای توییت هایی که **parentId** دارند، فیلد **parentContent** را حذف کنید.
4. پرتکرارترین و کم تکرارترین هشتگ را بیابید.
5. ده هشتگ (که معمولا همان نماد بورسی است) پر استفاده هر روز را بیابید. (بازه زمانی جزء ورودی های کوئری خواهد بود).
6. فعالترین کاربر هر روز را به همراه تعداد توییت های انجام شده، پیدا کنید.

خروجی گام چهارم

دستور نوشته شده، خروجی و زمان اجرا