



به نام خدا



دانشگاه تهران

دانشکده مهندسی برق و کامپیوتر

تحلیل ها و سیستم های داده های حجیم

پروژه ی پایانی

نام و نام خانوادگی	حسین قاسمی رامشه مسعود رحیمی سبحان راستی
شماره دانشجویی	810198224 810198161 810198160
تاریخ ارسال گزارش	29 تیر 1400

فهرست

4.....	مقدمه
5.....	گام اول - دریافت اطلاعات و Preprocess
5.....	مقدمه
5.....	بیشینه ی مثبت
7.....	معماری کافکا
10.....	گام دوم - persistence
10.....	گرفتن داده از کافکا
16.....	گام سوم - Cassandra
16.....	مقدمه
16.....	جدول اول
16.....	جدول دوم
17.....	جدول سوم
17.....	کد ها و توابع پایتون
22.....	گام چهارم - Redis
22.....	مقدمه
22.....	توضیحات کد پایتون پردازش داده ها
25.....	نمایش داده ها به وسیله ی Flask
27.....	گام پنجم - Clickhouse
27.....	مقدمه
27.....	توضیحات کد پایتون و ارتباط با Clickhouse
30.....	گام ششم - Apache Superset
30.....	مقدمه

30..... نمودارها و موارد خواسته شده

34..... پیوست - روند اجرای برنامه

در این پروژه با نحوه ی واکشی و پردازش داده های کلان و استفاده از پایگاه داده های مختلف برای پردازش، تحلیل و نمایش اطلاعات خروجی از آن ها آشنا شدیم. در هر قسمت چالش ها و مراحل کار از ورود داده تا نمایش نهایی را توضیح داده ایم. این [لینک](#) نیز به منظور دسترسی به کدهای این پروژه است.

گام اول - دریافت اطلاعات و Preprocess

مقدمه

برای داده‌های ورودی مشتاق بودم تا بتوانم داده‌های توییت‌ر را دریافت کرده و آن‌ها را پردازش کنم؛ پس از ایجاد کد جهت دریافت داده‌های توییت‌ر و استفاده از API‌های داخل دستورکار، ارور گرفتم و متوجه شدم که نیاز به دریافت API از توییت‌ر است. سپس درخواست API کردم. اما توییت‌ر موافقت نکرد و گفت فعلاً امکان‌پذیر نیست. سپس به سراغ ساخت بات برای فروش رفتم. پس از ایجاد حساب و ساخت بات برای فروش، توی گوگل سرچ کردم ولی موفق نشدم طریقه عضو کردن بات در کانال‌ها را پیدا کنم و انگلیسی هم نمی‌توانستم سرچ کنم.

سپس به سراغ تلگرام رفتم و چند روزی روی این موضوع سرچ می‌کردم. چندجایی در اینترنت خواندم که امکان عضو کردن بات در کانال‌هایی که خوم ادمین آن‌ها هستم وجود دارد و من نمی‌توانم بات را در کانال دیگران عضو کنم. سپس متوجه شدم که می‌توانم پیام‌های حساب کاربری شخصی خودم را به کمک دریافت api و کتابخانه telethon وارد پایتون کنم. به این شکل امکان دریافت پیام کانال‌های تلگرام وجود داشت. اما توی قسمتی از documentation نوشته بود که اگر سرعت دریافت پیام، علی‌الخصوص برای کشورهای حساس مثل ایران و روسیه، بالا باشد، امکان بسته شدن اکانت شما وجود دارد ([لینک مطلب](#) ذکر شده).

به همین دلیل ریسک نکردم و به سراغ توییت‌های سهام‌یاب رفتم.

بیشینه‌ی مثبت

ابتدا کافکا را نصب کرده و به کمک command‌های زیر server و zookeeper را اجرا می‌کنم.

```
#### 1st cmd in kafka dir
```

```
\.bin\windows\zookeeper-server-start.bat .\config\zookeeper.properties
```

```
#### 2nd cmd in kafka dir
```

```
\.bin\windows\kafka-server-start.bat .\config\server.properties
```

حال با فراخوانی مداوم API داده‌ها را از سهام‌یاب دانلود می‌کنیم و سپس آن‌هایی که id تکراری ندارند را به کمک producer وارد topic کافکا به نام raw_data می‌کنیم. در همین حین، با یک consumer داده‌های خام را از تاپیک raw_data دریافت کرده و روی هر توییت پیش پردازش‌های گفته شده را اعمال می‌کنیم.

با دریافت هر توییت، ابتدا یک timestamps و UUID برای آن‌ها تولید کردم. به کمک یک تابع هشتگ‌های موجود در content هر توییت را یافته و آن را به اطلاعات همان توییت اضافه می‌کنیم. همین‌طور لینک‌ها هم استخراج شد.

جهت محاسبه کلمات کلیدی، از معیار tf/idf استفاده شده است. ابتدا لیستی از ایست‌واژه‌ها یافتم و با حذف کردن آن‌ها از متن توییت، tf/idf را محاسبه کردم. در tf/idf فراوانی هر کلمه را در هر توییت می‌یابیم و همین‌طور فراوانی همان کلمه را در همه‌ی توییت‌ها بررسی می‌کنیم. کلمه‌ای از tf/idf بالا برخوردار است (کلمه کلیدی است) که در یک توییت وجود داشته باشد اما در بقیه توییت‌ها کمتر تکرار شده باشد. سپس برای هر توییت علاوه بر کلمات خاصی که صورت سوال ذکر کرده است، دو کلمه‌ای که tf/idf بزرگتری دارند را به عنوان کلمات کلیدی در نظر گرفتیم.

هر توییت ما یک دیکشنری است. در تصویر 1 key, value یک توییت را نمایش داده‌ام و در سمت راست تصویر اطلاعاتی که به عنوان متادیتا ذخیره خواهد شد، قرار دارد.

Key	Type	Size	Value
_meta	dict	7	{'UUID': '69e59ed4-6f14-47e3-9965-393b614015b2', 't...
content	str	1	#غنوش
finalPullDatePersian	str	1	
id	str	1	69e59ed4-6f14-47e3-9965-393b614015b2
links	list	0	[]
scoredPostDate	str	1	1624877945362
sendTime	str	1	2021-06-28T10:58:51Z
sendTimePersian	str	1	1400/04/07 15:28
senderName	str	1	Amir
senderProfileImage	str	1	06592890-4794-4b1d-adda-2979c8702fbc
senderUsername	str	1	amir5665
type	str	1	twit

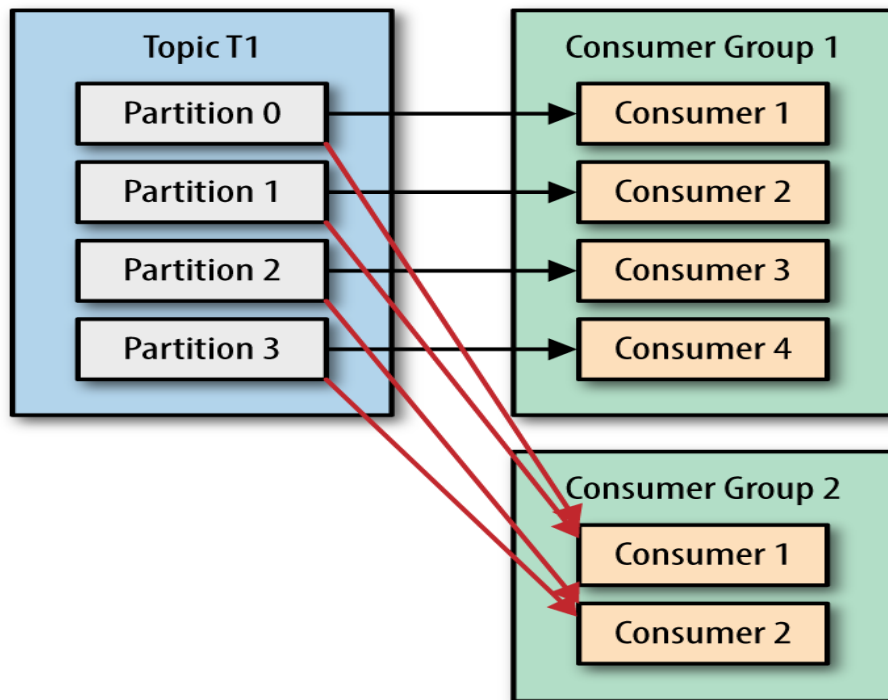
Key	Type	Size	Value
UUID	str	1	69e59ed4-6f14-47e3-9965-3...
hashtags	list	1	['غنوش']
key_words	list	2	['درش', 'خبری']
offset	int	1	5708
partition	int	1	0
timestamp	int	1	1626592694836
topic	str	1	test11

تصویر 1- کلید های داده ها

در تصویر 1 می‌بینیم که partition برابر صفر و offset برابر ۵۷۰۸ ایت. یعنی ایت توییت در تاپیک test11 و در پارتیشن اول آن و همین‌طور در خانه‌ی ۵۷۰۸ پارتیشن اول قرار دارد. در نهایت به کمک یک producer جدید داده‌های پیش پردازش شده را در تاپیک proccesed_data قرار می‌دهیم.

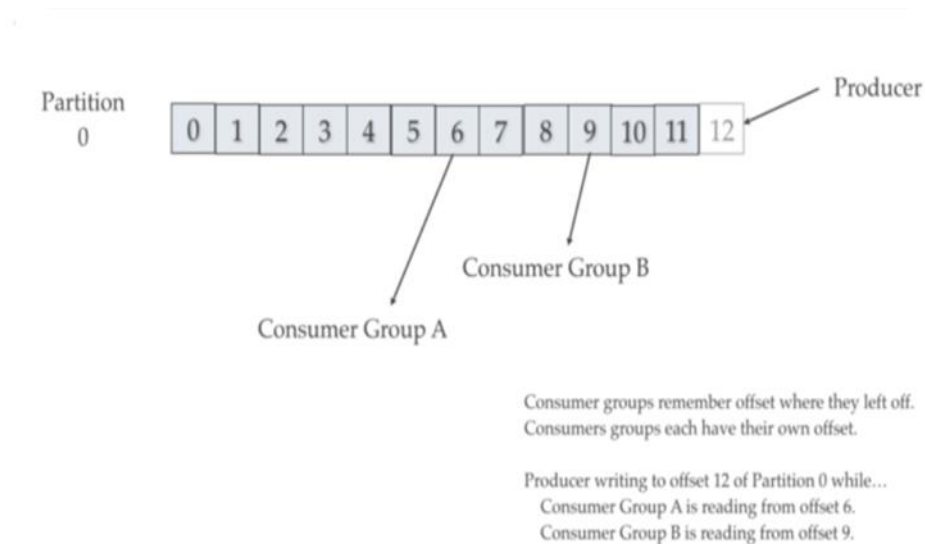
معماری کافکا

کافکا شامل جمعی از record, topic, consumer, producer, broker, log و cluster است و هر record دارای key (اختیاری)، value و timestamp است. record هایی که وارد کافکا می‌شوند immutable هستند؛ یعنی producer به تنهایی فقط می‌تواند پیام‌ها را به تاپیک topic اضافه کند و امکان تغییر در record های موجود را ندارد. هر تاپیک یک log دارد که محل ذخیره‌ی topic روی دیسک است. هر topic متشکل از تعدادی partition است. هر پارتیشن متشکل از تعدادی offset است (عکس زیر درک بهتری ایجاد می‌کند). هر record وارد topic مورد نظر و سپس وارد partition مورد نظر می‌شود و در آنجا به ترتیب درون offsetها ذخیره می‌شود. برای آنکه چندین consumer بتوانند همزمان record دریافت کنند، نیاز به وجود partition داریم. در عکس زیر نحوه‌ی اتصال consumer ها به partition را می‌بینیم. با ایجاد consumer با group id متفاوت، می‌توان همزمان از یک تاپیک پیام دریافت کرد.



تصویر 2- معماری کافکا

هر پارتیشن بدین شکل است:



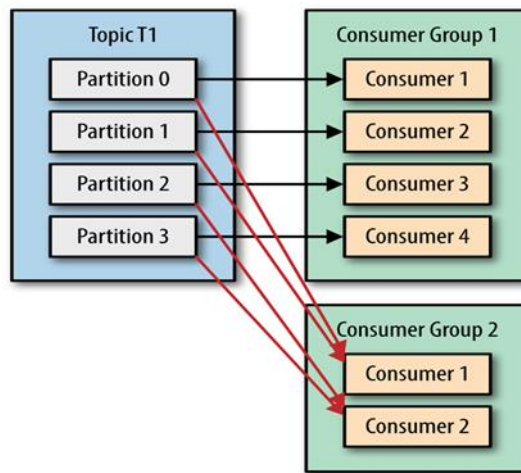
تصویر 3- پارتیشن بندی کافکا

Producer برای تولید جریان داده و ارسال آن به topic استفاده می‌سود. Consumer پیام‌ها را از partition کافکا می‌خواند (هر topic شامل تعدادی partition است).
Cluster کافکا متشکل از تعدادی broker است و هر broker شامل تعدادی topic است.

گام دوم – persistence

گرفتن داده از کافکا

در این قسمت داده‌های پیش پردازش شده‌ی مرحله قبل را به کمک consumer با group id معین دریافت می‌کنیم (چون در گام‌های بعدی از تایپیک processed_data داده دریافت خواهیم کرد و در واقع داریم به صورت همزمان از partitionهای واحد توسط چند consumer پیام دریافت می‌کنیم، تعیین group id حیاتی است). این موضوع به تصویر 4 اشاره دارد:

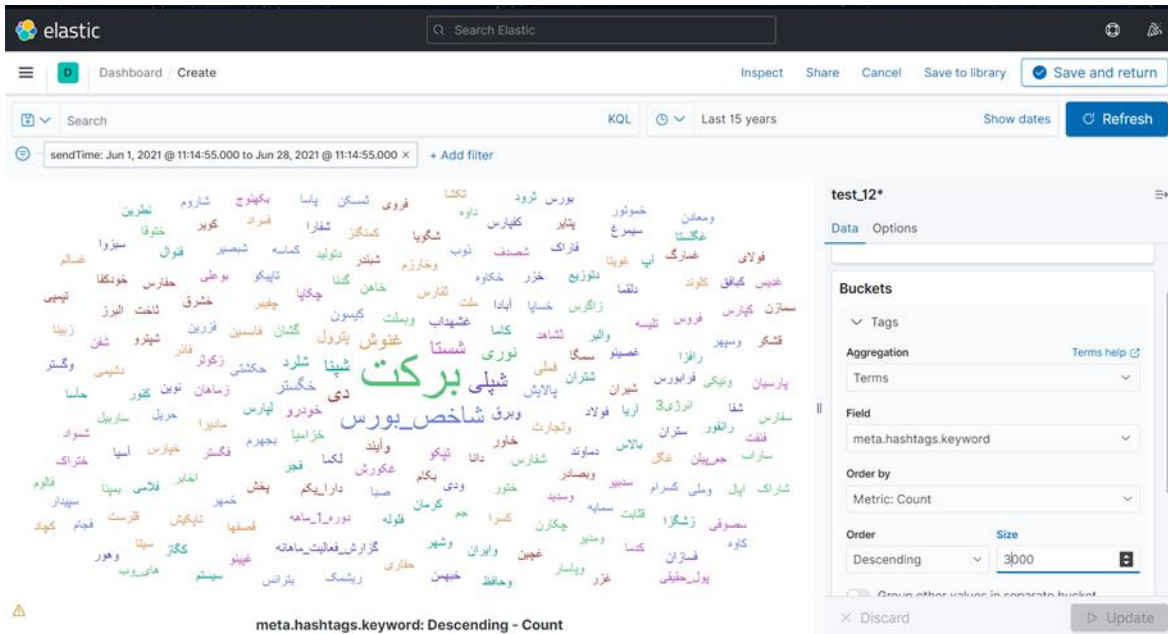


تصویر 4- group id

پس از دریافت داده ها از کافکا، داده‌ها را وارد index الاستیک سرچ می‌کنیم. سپس به کمک کیبانا، داشبوردهای خواسته شده را نمایش می‌دهیم.

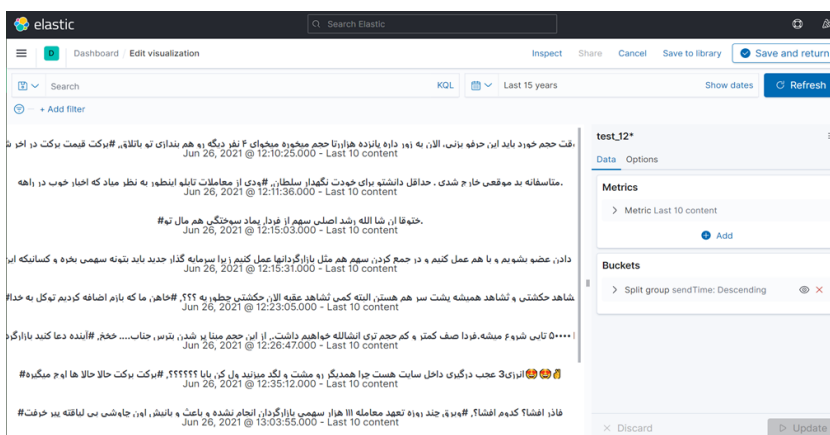
ابر کلمات یک کانال یا خبرگزاری خاص در یک بازه زمانی

از آنجایی که ما از داده‌های سهام‌یاب استفاده می‌کنیم، من ابر کلمات هشتگ‌ها در یک بازه‌ی خاص را رسم کرده‌ام.



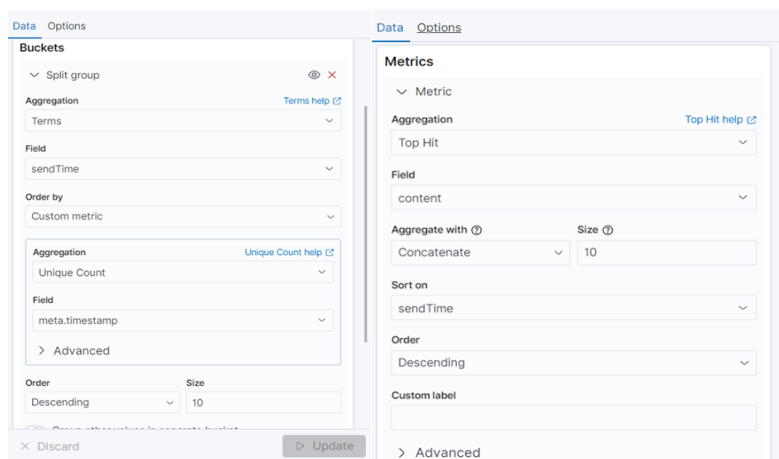
تصویر 5- ابر کلمات در الاستیک سرچ

متن ده پست اخیری که دریافت شده است



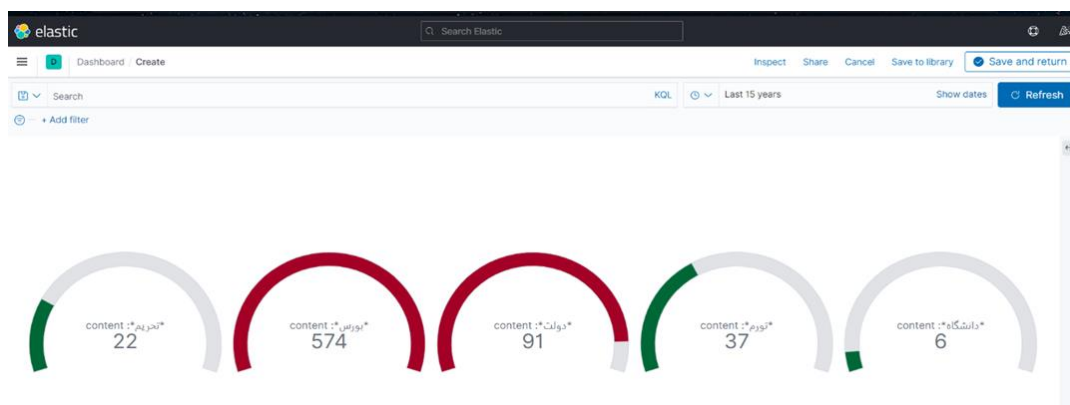
تصویر 6- متن ده پست اخیر در الاستیک سرچ

پارامترهای تعیین شده جهت نمایش داشبورد در تصاویر زیر آورده شده است (یعنی metrics و buckets).

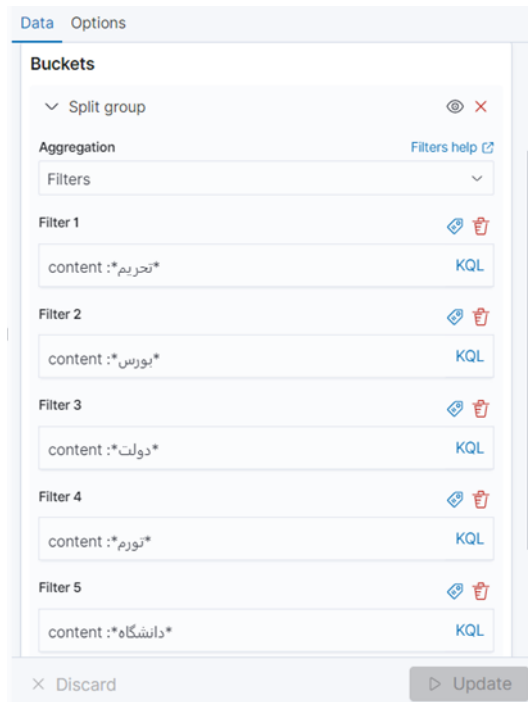


تصویر 7- پارامترهای داشبورد الاستیک سرچ

تعداد پست‌های ارسال شده به ازای چند تا از کلمات کلیدی خاص که در مرحله قبل مشخص شده است در یک بازه زمانی خاص برای این منظور نتیجه و bucket استفاده شده را نمایش دادیم.

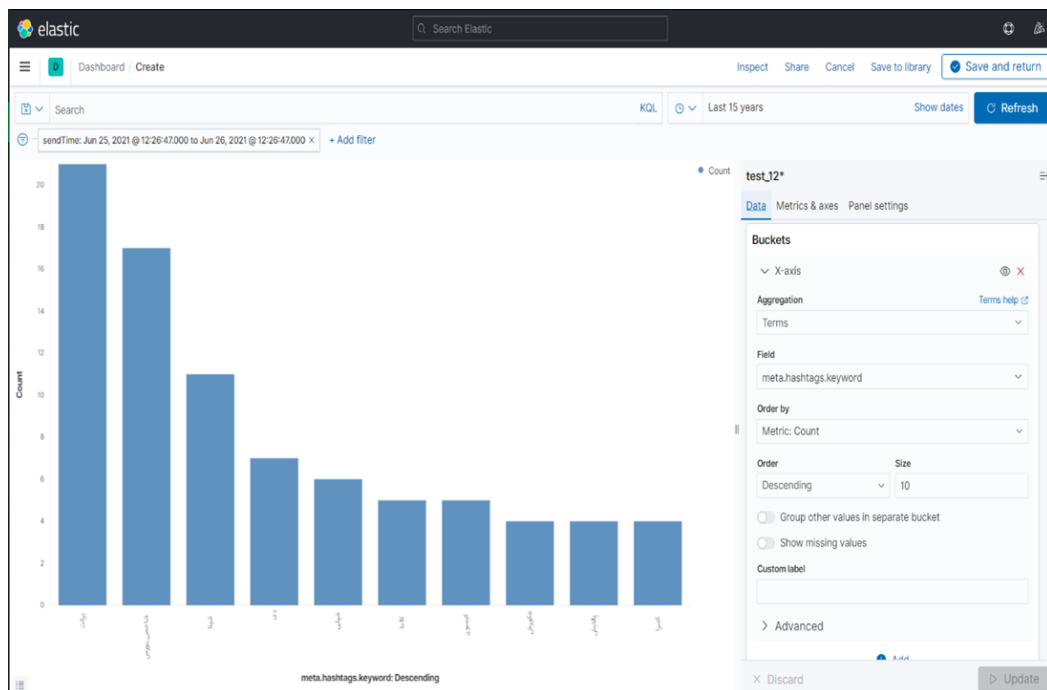


تصویر 8- آمار کلمات کلیدی در پست‌ها در الاستیک سرچ



تصویر 9- باکت برای آمار کلمات کلیدی در پست ها در الاستیک سرچ

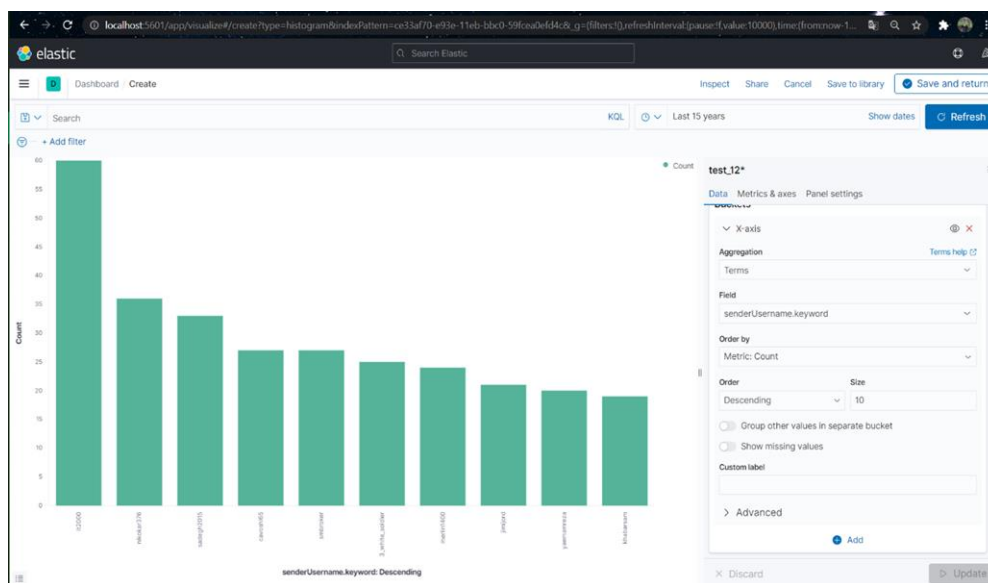
ده هشتگ بیشتر استفاده شده در پست‌های تمام کانال‌ها در یک بازه‌ی زمانی با تعداد تکرار هر هشتگ (یک نمودار ستونی) مثلاً هشتگ‌های بیشتر استفاده شده در یک روز اخیر



تصویر 10- نمودار ده هشتگ بیشتر استفاده شده در پست ها در الاستیک سرچ

یک نمودار به انتخاب خودتان

در این قسمت تعداد توییت‌های یک کاربر را در یک نمودار ستونی نمایش داده‌ام.



تصویر 10- نمودار تعداد توییت های کاربر در الاستیک سرچ

گر به دنبال تمام پست‌های حاوی یک کلمه خاص از یک خبرگزاری یا کانال خاص در یک بازه مشخص هستیم، چه دستوری باید بنویسیم (و یا یک هشتگ خاص یا یک کاربر خاص در توئیته‌ها)

```

# final proj
#1 GET /test12/_search
{
  "query": {
    "bool": {
      "must": [
        {
          "match": {
            "meta.hashtags": "#رکت"
          }
        }
      ]
    }
  },
  "range": {
    "sendTime": {
      "gte": "2020-01-01T00:00:00",
      "lte": "now"
    }
  }
}

#2 GET /test13/_search
{
  "query": {
    "bool": {
      "must": [
        {
          "match": {
            "content": "#رکت"
          }
        }
      ]
    }
  }
}

1 {
  "took": 11,
  "timed_out": false,
  "shards": {
    "total": 1,
    "successful": 1,
    "skipped": 0,
    "failed": 0
  },
  "hits": {
    "total": {
      "value": 515,
      "relation": "eq"
    },
    "max_score": 3.1629546,
    "hits": [
      {
        "_index": "test_12",
        "_type": "tweet",
        "_id": "Uql3w3oBHFru3qT7c_EC",
        "_score": 3.1629546,
        "_source": {
          "id": "1feb3bd5-8d75-46ab-a14c-dd16bfda3660",
          "sendTime": "2021-06-26T07:49:16Z",
          "sendTimePersian": "1400/04/05 12:19",
          "senderName": "Anonymous ",
          "senderUsername": "amirhoseinkh7",
          "senderProfileImage": "default",
          "content": ""#رکت""
        }
      }
    ]
  }
}

```

تصویر 11- دستور هشتگ خاص در توئیته‌ها در الاستیک سرچ

تعداد توئیته‌های ارسالی به ازای یک کلمه خاص را به ازای یک هشتگ خاص در توئیته‌ها در یک بازه زمانی

```

#2 GET /test12/_search
{
  "query": {
    "bool": {
      "must": [
        {
          "match": {
            "meta.hashtags": "#شخص بورس"
          }
        }
      ]
    }
  }
}

1 {
  "took": 4,
  "timed_out": false,
  "shards": {
    "total": 1,
    "successful": 1,
    "skipped": 0,
    "failed": 0
  },
  "hits": {
    "total": {
      "value": 8,
      "relation": "eq"
    },
    "max_score": 7.2275667,
    "hits": [
      {
        "_index": "test_12",
        "_type": "tweet",
        "_id": "Wall3w3oBHFru3qT73_dA",
        "_score": 7.2275667,
        "_source": {
          "id": "47a92453-e871-4485-9e37-b1ff2820c5c6",
          "sendTime": "2021-06-26T12:35:59Z",
          "sendTimePersian": "1400/04/05 17:05",
          "parentSendTime": "2021-04-19T18:51:42Z",
          "parentSendTimePersian": "1400/01/30 23:21",
          "parentId": "271539403",
          "parentSenderName": "سعید آقایی",
          "parentSenderUsername": "sfaghahi",
          "parentSenderProfileImage": "1ec38f47-70ac-4d3a-ae77"
        }
      }
    ]
  }
}

```

تصویر 12- تعداد توئیته‌ها به ازای هشتگ خاص در توئیته‌ها در الاستیک سرچ

گام سوم – Cassandra

مقدمه

بنده در این سوال سه جدول (البته استفاده از جدول درست نیست) ایجاد کردم یکی از آن‌ها بنابر سوالات داده شده طراحی شده است. در هر سه جدول کلید اول را کانال قرار داده‌ایم البته ما در این داده‌ها همه را از سهامیاب گرفته‌ایم پس فقط یک کانال داریم. بنده از پایتون برای اجرای کویری‌های کاساندرا استفاده کردم و این کار را با استفاده از کتابخانه cassandra انجام داده‌ام.

جدول اول

در جدول اولی چون می‌خواهیم به سوالات زمان ارسال پاسخ دهیم، پس sendTime را کلید بعدی قرار داده‌ایم:

```
session.execute("""
    CREATE TABLE IF NOT EXISTS fullData (
        canal text,
        sendTime timestamp,
        id text ,
        senderName text ,
        content text ,
        senderUsername text ,
        PRIMARY KEY (canal, sendTime, id) );
""")
```

تصویر 13- جدول اول کاساندرا

جدول دوم

در جدول دوم هشتگ را کلید دوم گذاشتیم و بعد زمان ارسال را به این دلیل که بتوانیم به راحتی برای هر هشتگ در یک بازه زمانی درخواست شده داده‌های مربوطه را برگردانیم.


```
session.execute("""
    CREATE TABLE IF NOT EXISTS hashtag (
        canal text,
        hashtag text,
        sendTime timestamp,
        id text ,
        senderName text ,
        content text ,
        senderUsername text ,
        PRIMARY KEY (canal, hashtag, sendTime, id) );
""")
```

تصویر 14- جدول دوم کاساندرا

جدول سوم

در جدول سوم کلمات کلیدی را کلید دوم گذاشتیم و بعد زمان ارسال را به این دلیل که بتوانیم به راحتی برای هر کلمه کلیدی در یک بازه زمانی درخواست شده داده‌های مربوطه را برگردانیم.

```
session.execute("""
    CREATE TABLE IF NOT EXISTS keyWord (
        canal text,
        keyword text,
        sendTime timestamp,
        id text ,
        senderName text ,
        content text ,
        senderUsername text ,
        PRIMARY KEY (canal, keyword, sendTime, id) );
""")
```

تصویر 15- جدول سوم کاساندرا

کد ها و توابع پایتون

در داده‌ای که وارد میشود ابتدا آن را به کاساندرا اضافه میکنیم:

```

insert_fullData = session.prepare("""
    INSERT INTO fullData (
        canal,
        sendTime,
        id,
        senderName,
        content,
        senderUsername
    )
    VALUES (?, ?, ?, ?, ?, ?)
    IF NOT EXISTS;
""")

insert_hashtag = session.prepare("""
    INSERT INTO hashtag (
        canal,
        hashtag,
        sendTime,
        id,
        senderName,
        content,
        senderUsername
    )
    VALUES (?, ?, ?, ?, ?, ?, ?)
    IF NOT EXISTS;
""")

insert_keyWord = session.prepare("""
    INSERT INTO keyWord (
        canal,
        keyword,
        sendTime,
        id,
        senderName,
        content,
        senderUsername
    )
    VALUES (?, ?, ?, ?, ?, ?, ?)
    IF NOT EXISTS;
""")

```

تصویر 16- کد ورود داده کاساندرا

برای استفاده از داده‌های وارد شده بنده شش تابع نوشته‌ام که سه عدد از آن‌ها دو بازه زمانی کار میکنند که هر کدام از این سه با یک جدول کار میکند. سه تابع بعدی سوالات خواسته شده را پشتیبانی میکند به این صورت که دو ورودی دارد که یکی تعداد (ساعت، روز، هفته، ماه، سال) و دومی اگر نیاز بود هشتگ یا کلمه کلیدی مربوطه میباشد.

ابتدا کدهایی را مشاهده میفرمایید که بازه زمانی را میگیرند:

```
def all_post(timestamp1 , timestamp2, canal = 'sahamyab' ):
    cluster = Cluster(['127.0.0.1'])
    session = cluster.connect()
    session.set_keyspace(KEYSPACE)
    code = session.prepare("""
        SELECT * FROM fullData
        WHERE canal = ? AND sendTime >= ? AND sendTime <= ?;
    """)
    result = session.execute(code, [canal, timestamp1, timestamp2])
    return result

def all_post_by_hashtag(hashtag, timestamp1 , timestamp2, canal = 'sahamyab'):
    cluster = Cluster(['127.0.0.1'])
    session = cluster.connect()
    session.set_keyspace(KEYSPACE)
    code = session.prepare("""
        SELECT * FROM hashtag
        WHERE canal = ? AND hashtag = ? AND sendTime >= ? AND sendTime <= ?;
    """)
    result = session.execute(code, [canal, hashtag, timestamp1, timestamp2])
    return result

def all_post_by_keyWord(keyword, timestamp1 , timestamp2, canal = 'sahamyab' ):
    cluster = Cluster(['127.0.0.1'])
    session = cluster.connect()
    session.set_keyspace(KEYSPACE)
    code = session.prepare("""
        SELECT * FROM keyWord
        WHERE canal = ? AND keyWord = ? AND sendTime >= ? AND sendTime <= ?;
    """)
    result = session.execute(code, [canal, keyword, timestamp1, timestamp2])
    return result
```

تصویر 17- کد بازه زمانی کاساندرا

در ادامه کدهای داده‌هایی را مشاهده میکنید که تعداد میگیرد و بنا بر تعداد برای ساعت، ماه و را محاسبه میکند و بر میگرداند و در کدها پیشبینی بزرگ شدن عدد ماه یعنی بزرگتر از 12 یا کوچکتر از 1 شده است:

```
def get_post(n):
    timestamp2 = datetime.today()
    timestamp1 = timestamp2 - timedelta(hours=n)
    n_hour = all_post(timestamp1 , timestamp2)

    timestamp1 = timestamp2 - timedelta(days=n)
    n_day = all_post(timestamp1 , timestamp2)

    timestamp1 = timestamp2 - timedelta(weeks=n)
    n_week = all_post(timestamp1 , timestamp2)

    m = timestamp2.month - int(n % 12)
    y = timestamp2.year - int(n / 12)
    if m <= 0:
        m = 12 + m
        y -= 1
    timestamp1 = datetime.strptime(f'{y}-{m}-{timestamp2.day}', '%Y-%m-%d')
    n_month = all_post(timestamp1 , timestamp2)

    timestamp1 = datetime.strptime(f'{timestamp2.year-n}-{timestamp2.month}-{timestamp2.day}', '%Y-%m-%d')
    n_year = all_post(timestamp1 , timestamp2)

    return n_hour, n_day, n_week, n_month, n_year
```

تصویر 18- کد محاسبه ی تعداد برای زمان مختلف کاساندر بخش اول

```

def get_post_by_hashtag(n, hashtag):
    timestamp2 = datetime.today()
    timestamp1 = timestamp2 - timedelta(hours=n)
    n_hour = all_post_by_hashtag(hashtag, timestamp1 , timestamp2)

    timestamp1 = timestamp2 - timedelta(days=n)
    n_day = all_post_by_hashtag(hashtag, timestamp1 , timestamp2)

    timestamp1 = timestamp2 - timedelta(weeks=n)
    n_week = all_post_by_hashtag(hashtag, timestamp1 , timestamp2)

    m = timestamp2.month - int(n % 12)
    y = timestamp2.year - int(n / 12)
    if m <= 0:
        m = 12 + m
        y -= 1
    timestamp1 = datetime.strptime(f'{y}-{m}-{timestamp2.day}', '%Y-%m-%d')
    n_month = all_post_by_hashtag(hashtag, timestamp1 , timestamp2)

    timestamp1 = datetime.strptime(f'{timestamp2.year-n}-{timestamp2.month}-{timestamp2.day}', '%Y-%m-%d')
    n_year = all_post_by_hashtag(hashtag, timestamp1 , timestamp2)

    return n_hour, n_day, n_week, n_month, n_year

```

تصویر 19- کد محاسبه ی تعداد برای زمان مختلف کاساندرای بخش دوم

```

def get_post_by_keyWord(n, keyword):
    timestamp2 = datetime.today()
    timestamp1 = timestamp2 - timedelta(hours=n)
    n_hour = all_post_by_keyWord(keyword, timestamp1 , timestamp2)

    timestamp1 = timestamp2 - timedelta(days=n)
    n_day = all_post_by_keyWord(keyword, timestamp1 , timestamp2)

    timestamp1 = timestamp2 - timedelta(weeks=n)
    n_week = all_post_by_keyWord(keyword, timestamp1 , timestamp2)

    m = timestamp2.month - int(n % 12)
    y = timestamp2.year - int(n / 12)
    if m <= 0:
        m = 12 + m
        y -= 1
    timestamp1 = datetime.strptime(f'{y}-{m}-{timestamp2.day}', '%Y-%m-%d')
    n_month = all_post_by_keyWord(keyword, timestamp1 , timestamp2)

    timestamp1 = datetime.strptime(f'{timestamp2.year-n}-{timestamp2.month}-{timestamp2.day}', '%Y-%m-%d')
    n_year = all_post_by_keyWord(keyword, timestamp1 , timestamp2)

    return n_hour, n_day, n_week, n_month, n_year

```

تصویر 20- کد محاسبه ی تعداد برای زمان مختلف کاساندرای بخش سوم

مقدمه

در این قسمت پس از نصب و راه اندازی Redis در حالیکه سرور و کلاینت Redis در حال اجرا است به وسیله ی کدهای پایتون به Redis متصل شده و داده های مورد نیاز را به Redis میفرستیم. برای دریافت آنلاین داده ها از کافکا به سراغ فایلی میرویم که همواره در حال آپدیت شدن است و داده های به روز را در فایلی در آدرس کد مد نظر ما ذخیره می کند. به منظور این که حجم این فایل از حدی بیشتر نشود بایستی صرفا داده های جدید را در فایل ذخیره کرد. در انتها و پس از انجام پردازش های مورد نیاز و گرفتن نتایج آن ها را از طریق صفحه ی HTML ای که توسط فلسک آماده شده به صورت مداوم نمایش میدهیم.

توضیحات کد پایتون پردازش داده ها

برای اینکه تمامی مراحل به صورت همروند انجام شده و نمایش داده بر روی وب اپلیکیشن هم به روز باشد از یک Thread برای گرفتن داده، ارتباط با Redis و بدست آوردن نتایج استفاده میکنیم و نمایش بر روی اپلیکیشن نیز توسط Thread اصلی انجام میگردد.

Thread پردازش تابعی به اسم process که در تصاویر زیر مشخص است را اجرا میکند. در این تابع به صورت مداوم و درون حلقه کارهای مورد نیاز انجام میگردد. ضمنا Thread طوری طراحی شده که با تمام شدن کار در Thread اصلی آن نیز به پایان برسد. در تابع process ابتدا داده ها را از فایل مربوطه خوانده و نگه میداریم. سپس به ازای تمام توییت های دریافتی کارهای پردازشی و اضافه کردن آن ها به Redis را انجام میدهیم. برای این منظور کلیدی منحصر به آن توییت خاص و مرکب از اسم فرستنده ی توییت و زمان ارسال آن تهیه کرده و سپس در صورتی که آن کلید در Redis موجود نبود، کلید مدنظر را با یک عدد یک برای شمارش و همچنین زمان انقضای یک هفته ای به Redis اضافه میکنیم و در غیر اینصورت تعداد آن کلید خاص را تنها اضافه میکنیم. همین کار را برای کلید متشکل از زمان توییت هم انجام میدهیم. در صورتی که لیست مربوط به هشتگ های توییت خالی نبود این لیست را به همراه یک کلید مجزا ذخیره مینماییم. همین کار را برای ساخت لیست هزارتایی هشتگ ها هم انجام داده و هربار بعد از push کردن داده جدید اگر طول لیست از هزار بالاتر شد آن را Trim کرده و به اندازه ی 1000 برمیگردانیم. بعد از ذخیره سازی هشتگ ها خود متن توییت ها را نیز در لیستی دیگر اما این بار با حداکثر اندازه ی 100 ذخیره میکنیم.

برای گرفتن نتایج برای سوال اول که همان مجموع توییت های دریافتی در شش ساعت اخیر است، به ازای هر ساعت، کلید ساخته شده برای شمارش این قسمت را به وسیله ی کم کردن آن ساعت از زمان فعلی بازیابی کرده و تعداد توییت موجود در Redis با آن کلید را به دست می آوریم و حاصل را جمع میکنیم.

برای سوال دوم همان کار قبلی را اینبار به ازای مقادیر زمان های خاص در یک بازه ی زمانی دلخواه مثلاً یک روز گذشته تکرار میکنیم. برای سوال سوم همین کار اینبار برای هشتگ ها و برای بازه ی یک ساعت گذشته اجرا کرده و تعداد هشتگ ها را بدست می آوریم. با گرفتن موارد ذخیره شده با کلید لیست هشتگ ها، به هدف سوال چهار یعنی بدست آوردن یک لیست هزارتایی از هشتگ ها میرسیم. در نهایت به همین روش اما با کلیدی متفاوت لیست توییت ها را دریافت میکنیم. در انتها همه ی جواب ها را به همراه متن مناسب برای استفاده ی فلسک در متغیری global ذخیره میکنیم. در زیر تصاویر مربوط به کد این تابع را میبینید:

```
34 def process():
35     global Q1
36     global Q2
37     global Q3
38     global Q4
39     global Q5
40     global hashtagsList
41     global tweetsList
42     while True:
43         time.sleep(0.01)
44         f2 = open('preprocessed_data.json', encoding="utf8")
45         data = json.load(f2)
46         f2.close()
47
48         expirationMinutes = find_expiration_minutes(7)
49
50         for item in data:
51             # print(item)
52             key = "username: " + item["senderUsername"] + " in: " + find_date(item["sendTime"])
53             if redis.exists(key) == 0:
54                 redis.set(key, 1, ex=expirationMinutes)
55             else:
56                 redis.incr(key)
57
58             key = "tweets: " + find_date(item["sendTime"])
59             if redis.exists(key) == 0:
60                 redis.set(key, 1, ex=expirationMinutes)
61             else:
62                 redis.incr(key)
63
64             if "" not in item["_meta"]["hashtags"] and len(item["_meta"]["hashtags"]) != 0:
65                 key = "unique_hashtags: " + find_date(item["sendTime"])
66                 if redis.exists(key) == 0:
67                     redis.sadd(key, *item["_meta"]["hashtags"])
68                     redis.expire(key, expirationMinutes)
69                 else:
70                     redis.sadd(key, *item["_meta"]["hashtags"])
71
72             if redis.exists("hashtags_list") == 0:
73                 redis.lpush("hashtags_list", *item["_meta"]["hashtags"])
```

تصویر 21- بخش اول کد تابع process برای Redis

```

74         redis.expire("hashtags_list", expirationMinutes)
75     else:
76         redis.lpush("hashtags_list", *item["meta"]["hashtags"])
77         if redis.llen("hashtags_list") > 1000:
78             redis.ltrim("hashtags_list", 0, 999)
79
80     if redis.exists("tweets_list") == 0:
81         redis.lpush("tweets_list", item["content"])
82         redis.expire("tweets_list", expirationMinutes)
83     else:
84         redis.lpush("tweets_list", item["content"])
85         if redis.llen("tweets_list") > 100:
86             redis.ltrim("tweets_list", 0, 99)
87
88
89     lastHours = 6
90     now = datetime(2021, 7, 20, 23)
91
92     totalTweets = 0
93     for h in range(lastHours + 1):
94         goalDate = now - timedelta(hours=h)
95         key = "tweets: " + find_date_from_obj(goalDate)
96         if redis.exists(key) > 0:
97             totalTweets += int(redis.get(key).decode("utf-8"))
98
99     Q1 = "توییت درهالت شده است " + str(totalTweets) + " ساعت اخیر در "
100
101     dateFrom = datetime(2021, 7, 1, 18)
102     dateTo = datetime(2021, 7, 20, 23)
103
104     delta = timedelta(hours=1)
105     i = dateFrom
106     totalTweets = 0
107     while i <= dateTo:
108         key = "tweets: " + find_date_from_obj(i)
109         if redis.exists(key) > 0:
110             totalTweets += int(redis.get(key).decode("utf-8"))
111         i += delta

```

تصویر 22- بخش دوم کد تابع process برای Redis

```

112
113     Q2 = "درهالت شده است " + str(totalTweets) + " ساعت " + str(dateFrom) + " تا " + str(dateTo) + " از تاریخ "
114
115     lastHours = 1
116     totalHashtags = 0
117
118     for h in range(lastHours + 1):
119         goalDate = now - timedelta(hours=h)
120         key = "unique hashtags: " + find_date_from_obj(goalDate)
121         if redis.exists(key) > 0:
122             totalHashtags += redis.scard(key)
123
124     Q3 = "هشتگ درهالت شده است " + str(totalHashtags) + " ساعت اخیر در "
125
126
127     indexFrom = 0
128     indexTo = -1
129
130     #hashtagsList = []
131     for item in redis.lrange("hashtags_list", indexFrom, indexTo):
132         hashtagsList.append(item.decode("utf-8"))
133
134     Q4 = str(hashtagsList)
135
136
137
138     indexFrom = 0
139     indexTo = -1
140     # tweetsList = []
141     # print(tweetsList)
142     for item in redis.lrange("tweets_list", indexFrom, indexTo):
143         tweetsList.append(item.decode("utf-8"))
144
145     Q5 = "مجموعه توییت آخر این ها هستند "
146     # print(tweetsList)
147
148
149

```

تصویر 23- بخش سوم کد تابع process برای Redis

نمایش داده ها به وسیله ی Flask

در این بخش برای نمایش نتایج از Flask کمک گرفته و نتایج ذخیره شده در متغیرهای global را با تگ های مناسب HTML ترکیب کرده و آن را برای نمایش به Flask داده و نتیجه را اجرا مینماییم. در زیر کد مربوط به این بخش و نمونه های، از نتایج قابل مشاهده اند.

```

150 @app.route('/')
151 @app.route('/home')
152
153 def home():
154     s = ""
155     s = "<html><head><style>p.ex1 { border: 1px solid green;padding: 20px;margin: 20px;}</style></head><body>"
156     s = s + "<h2>" + Q1 + "</h2>"
157     s = s + "<h2>" + Q2 + "</h2>"
158     s = s + "<h2>" + Q3 + "</h2>"
159     s = s + "<h2>" + "مبارك الله ما حسنة" + "</h2>"
160     s = s + "<p class='ex1', style='background-color: #ffffff'>" + Q4 + "</p>"
161     s = s + "<h1>" + Q5 + "</h1>"
162     for x in tweetsList:
163         s = s + "<p class='ex1', style='background-color: #ffffff'>" + str(x) + "</p>"
164     s = s + "</tr>"
165     return "<html dir='rtl' Lang='ar'><body style='background-color: lightgrey'>" + s + "</body></html>"
166
167
168 if __name__ == '__main__': # Script executed directly?
169     process_thread = threading.Thread(target=process)
170     process_thread.setDaemon(True)
171     process_thread.start()
172     app.run(debug = True) # Launch built-in web server and run this Flask webapp
173

```

تصویر 24- کد تابع Flask و اجرای برنامه برای Redis



تصویر 25- بخش اول نمونه ای از خروجی سوال Redis

بخش: 'امیر خ'، اودی، 'شخص بورس'

صد توییت آخر این ها هستند:

تشفیر 800 تومان

بیشتر د نظرین سهامداران مطمئن بالتید شرکت به این روز انداخته چه خانواده هایی رو که بنیغت نکرده

دوست عزیز در این که با منفی باز میشه نشکی نیست حالا چه در بازار و بکل و بلبل و چه بازار فرسز. فقط باید امیدوار بود که با حداقل منفی باز شه و بتونه همون روز خوشو به صفت خرید برسونه. من چهارشنبه گفتیم که یکشنبه باز میشه ولی به حده و حد بازگشایی امروز رو داین.

بندی سالم طبق خبر های واصله جدید منبر. عامل بانک دی اعلام کرد شب گذشته کلید در سهام بانک دی رو از حسن کلید سار در بافت کرده که بعداز چند بار امتحان جهت سوهان کاری کلید به مرکز امور کلید بانک مرکزی مودت شده که انشاء الله طی چند روز آینده جهت باز شدن سهام دی آماده و ارسال میشود سهام داران محترم هیچگونه حکس العملی نشان نداده و در صورت فشار های مختلف اقتصادی فقط به خداوند متعال توکل کنند

حالا هزارتا سیگنال دادم ... بیست تا ترم رده شده باشه میشه دو دهم در صد.... برگردی اون ۱۵ تا دیگه هم پیدا می کنی ... خلیج مهم ابوبرقی هست که حسابش داره میشه ۱۵۰ در صد ... تو صورت کی تارکه خدا می دونه!!!!

بسیستم همچنان بر قدرت بر مدار صعود، خیر افزایش سرمایه بیاد ، مریم برای عبور از ۳۰۰ تومان ☺

بیشتر آن احتمال داره شاخص اول وقت فردا به پولیک به خط روند نزولیش بزنه ، اما بازم برمیگرده به بالا

تصویر 26- بخش دوم نمونه ای از خروجی سوال Redis

مقدمه

در این بخش چالش زیادی برای راه اندازی Clickhouse داشتیم، چرا که برای کار با این دیتاست نیاز به سیستم عامل لینوکس بود. البته ابزارهای دیگری در ویندوز این امکان را فراهم میکرد که خود آن ها نیز چالش های خاص خود را داشت. به هر صورت پس از نصب و راه اندازی Clickhouse در حالیکه سرور آن در حال اجراست به وسیله ی کد پایتون داده ها را در آن وارد کرده و جداول را ساختیم.

توضیحات کد پایتون و ارتباط با Clickhouse

در این بخش همانطور که در تصاویر زیر مشخص است پس از اتصال به سرور Clickhouse داده های پردازش شده و ریخته شده در فایل توسط کافکا به صورت آنلاین و مداوم را خوانده و پس از انتخاب بخش های مورد نیاز آن، حاصل را به Clickhouse اضافه مینماییم.

```
1 import json
2 from clickhouse_driver import Client
3 from datetime import datetime
4
5 # client = Client('localhost')
6 client = Client('localhost',
7               password='****@')
8
9 client.execute('DROP TABLE IF EXISTS default.sahamyab')
10 client.execute("""CREATE TABLE IF NOT EXISTS default.sahamyab
11 (
12     id String,
13     sendTime DateTime,
14     sendTimePersian String,
15     hashtags Array(Nullable(String)),
16     key_words Array(Nullable(String)),
17     senderUsername String,
18     senderName String,
19     content String
20 )
21 ENGINE = MergeTree PARTITION BY toYYYYMMDD(sendTime) ORDER BY toYYYYMMDD(sendTime)
22 """)
23
24
25 def insert(jData):
26     data = [
27         jData['id'],
28         datetime.strptime(jData['sendTime'], '%Y-%m-%dT%H:%M:%SZ'),
29         jData['sendTimePersian'],
30         jData['_meta']['hashtags'],
31         jData['_meta']['key_words'],
32         jData['senderUsername'],
33         jData['senderName'],
34         jData['content']
35     ]
```

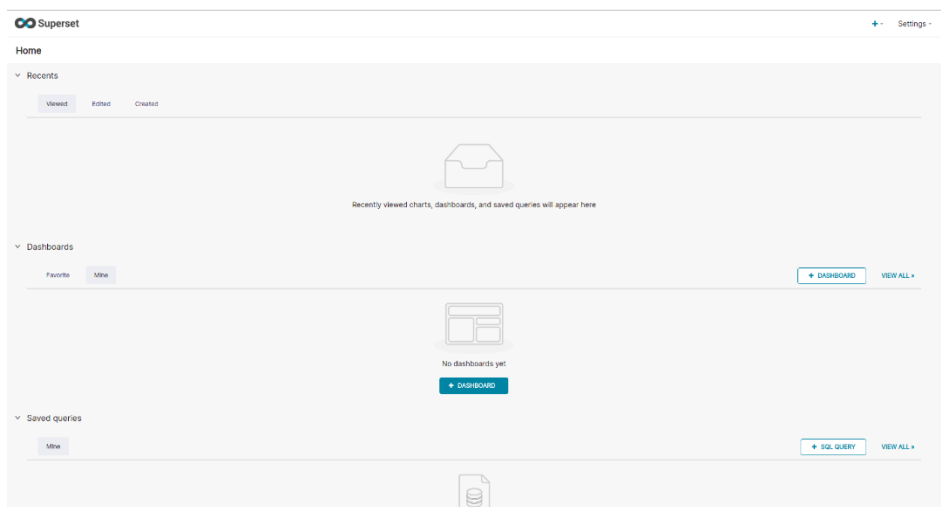
تصویر 27- بخش اول کد اضافه کردن داده به Clickhouse

مقدمه

در این بخش پس از اضافه شدن داده ها به Clickhouse و نصب و راه اندازی Apache Superset، Superset را به Clickhouse متصل کردیم. این اتصال در صورت داشتن رمز برای Clickhouse در ویندوز برای ما چالش برانگیز بود، چرا که نمیتوانستیم اتصال را برقرار کرده و از طرف دیگر نتوانستیم رمز را حذف نماییم. در نهایت در سیستم عامل لینوکس این مهم میسر شد.

نمودارها و موارد خواسته شده

با توجه به موارد خواسته شده در هر بخش کوئری مربوطه را داده و داده هارا از Clickhouse دریافت کرده سپس بعد از explore کردن داده ها و ذخیره ی آنها با انتخاب پارامتر ها و نمودارهای مناسب نتیجه را بدست آوردیم برای مواردی نیز به وسیله ی رابط کاربری superset و محیط گرافیکی آن و بدون نیاز به کوئری نمودار های خواسته شده را بدست آوردیم. همانطور که در دستور کار خواسته شده است سه داشبورد از موارد خواسته شده تهیه کردیم که آنها را در ادامه می بینیم.



تصویر 30- محیط superset

در تصاویر 31 و 32 جداول بدست آمده پس از در خواست های انجام شده را می بینید. در تصویر 31 جدول و نمودار نزولی برای بیشترین تکرار در بین هشتگ ها و همین مورد را در تصویر 32 برای کلمات کلیدی می بینید.

Superset Data - Charts Dashboards SQL Lab -

Q1-1

641 rows 00:00:00.57

hashtags

hashtags	COU
["u0628u0631u06a9u062a"]	1345
["u0634u0627u062au0635_u0628u0648u0631u0633"]	512
["u0634u067u0644u06cc"]	182
["u062fu06cc"]	149
["u0634u0633u062au0627"]	115
["u0634u0648u0648u0634"]	111
["u0634u067u0648u0627"]	102
["u0646u0648u0631u06cc"]	84
["u062au064fu0633u062au0631"]	71
["u0648u0628u06ccu0648u0627"]	55
["u0648u0628u0631u0642"]	48
["u0648u0627u062au0635_u0628u0648u0631u0633"]	46
["u067au062au0631u0648u0644"]	41
["u0634u0644u0631u0627"]	39
["u062au0648u062fu0631u0648"]	38
["u062au0632u0627u0645u06ccu0627"]	37

▼ Data

VIEW RESULTS VIEW SAMPLES

641 rows retrieved

hashtags

hashtags	COU
["u0628u0631u06a9u062a"]	1345
["u0634u0627u062au0635_u0628u0648u0631u0633"]	512
["u0634u067u0644u06cc"]	182
["u062fu06cc"]	149
["u062fu06cc"]	115

تصویر 31- هشتگ های دریافتی به همراه تعداد و نمودار آنها به صورت اسکی در محیط superset

Superset Data - Charts Dashboards SQL Lab -

Q1-2

1k rows 00:00:02.39

key words

key words	COU
["010"]	27
["u0628u0631u06a9u062a_u0645u0634u06a9u06cc"]	24
["u062au0627u06ccu06ccu0627_u0634u0627u062fu0633"]	10
["u0628u0648u0631u0633_u0645u0642u0627u06ccu0633u0647_u062fu0631u0635u062fu06cc"]	9
["u0627u0641u0632u0627u06ccu0634_u0633u0631u0649u0627u06ccu0647"]	8
["u0631u06ccu0627u0644_u0645u0627u0647u0647"]	7
["u0641u0631u0648u0634_u0645u0627u0647"]	7
["u0645u0628u0644u0634_u0631u06ccu0627u0644"]	6
["u0627u0644u0627u0631_u0627u0645u0627u0631u0627u062a_u0627u06ccu0631u0627u0648"]	6
["u0688u06a9u0627u0631u0646_u010"]	6
["u0637u0644u0627_u0648u0633u0628u062a_u0627u0648u0631u0647"]	6
["u062au0631u0648u062cu06cc_u0627u0628u0631"]	5
["u0628u0648u0631u0633_u0648u0645u0627u062fu0647u0627_u0641u0631u0648u0634"]	5
["u0631u0648u0648u0645u0627u06ccu06cc_u0648u0627u06a9u0633u0648"]	4
["u062cu0627u06ccu0632u0647_u0627u0636u0627u0641u0647"]	4
["u0628u0648u0631u0633_u0645u0627u0647_u062au0633u0647u06ccu0644u0627u062a"]	4

▼ Data

VIEW RESULTS VIEW SAMPLES

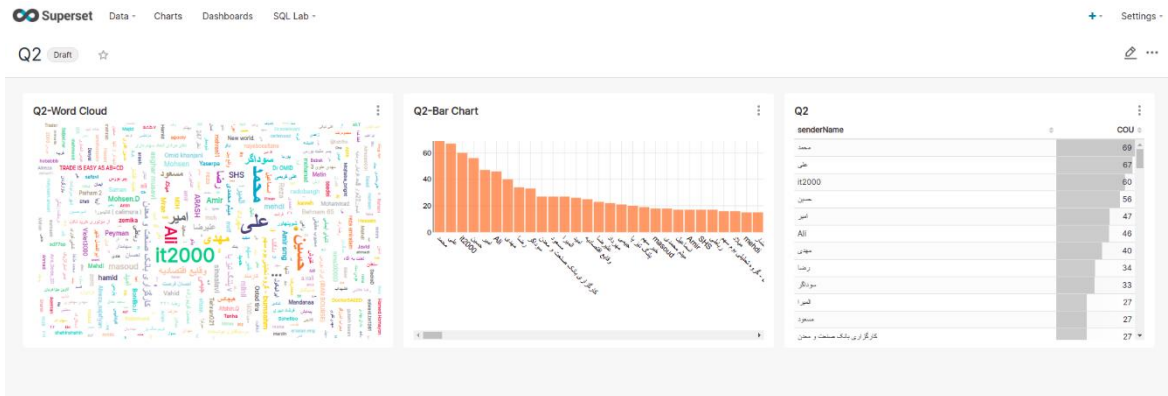
1k rows retrieved

key words

key words	COU
["010"]	27
["u0628u0631u06a9u062a_u0645u0634u06a9u06cc"]	24
["u062au0627u06ccu06ccu0621_u0634u0627u062fu0635"]	10
["u0628u0648u0631u0633_u0645u0642u0627u06ccu0633u0647_u062fu0631u0635u062fu06cc"]	9
["u0627u0641u0632u0627u06ccu0634_u0633u0631u0649u0627u06ccu0647"]	8

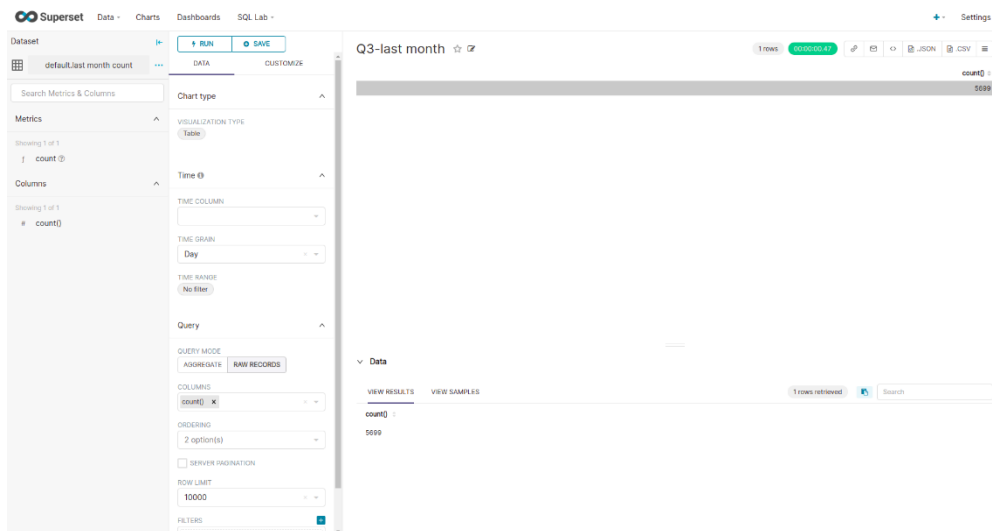
تصویر 32- کلمات کلیدی به همراه تعداد و نمودار آنها به صورت اسکی در محیط superset

در این بخش تعداد توییت های ارسال شده توسط هر کاربر را به دست آورده و به صورت نزولی مرتب کردیم و نمودار های مربوطه را کشیدیم. در ابر کلمات که یک ابزار کارآمد برای تحلیل داده های متنی است مشاهده میکنیم که کاربرانی با نام علی، محمد، سوداگر، it2000 از فعال ترین کاربرها در این حوزه بوده اند.

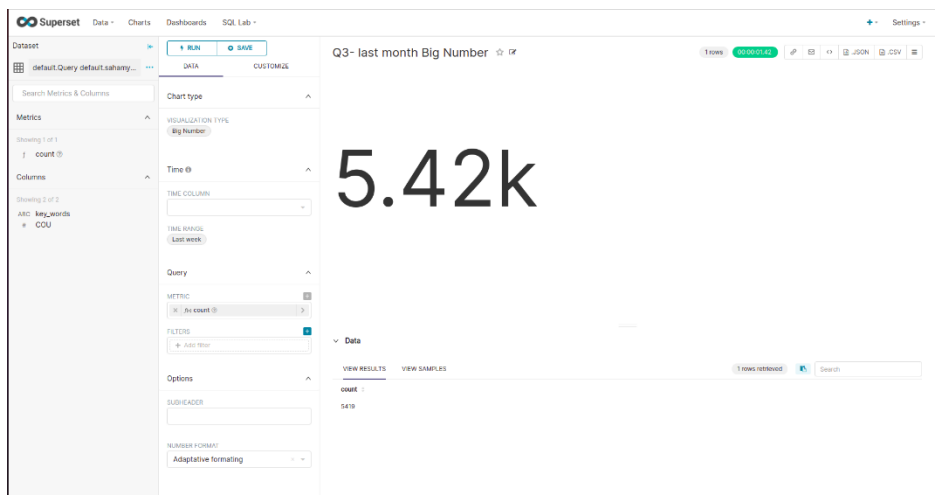


تصویر 33- تعداد توییت های کاربران به صورت نزولی در قالب جدول نمودار و ابرکلمات در محیط superset

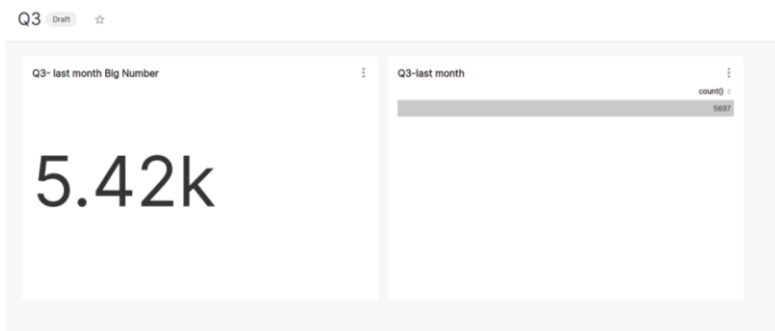
در ادامه تعداد توییت ها را در زمان ها مختلف بررسی کردیم که در تصویر 33 تعداد به دست آمده برای ماه اخیر را میبینیم و در تصاویر بعدی فرم این داده را به صورت Big Number کشیدیم.



تصویر 34- تعداد توییت ها در ماه اخیر در محیط superset



تصویر 35- Big Number تعداد توییت ها در ماه اخیر در محیط superset



تصویر 35- Big Number و تعداد توییت ها در ماه اخیر در محیط superset

پیوست - روند اجرای برنامه

به ازای هر قسمت کد مربوطه در پیوست آمده است.