



به نام خدا



دانشگاه تهران
دانشکده مهندسی برق و کامپیوتر
تحلیل داده‌های حجیم

پروژه پایانی

نام و نام خانوادگی	سبحان راستی
شماره دانشجویی	810198160
تاریخ ارسال گزارش	1400/04/29

فهرست گزارش سوالات

3.....	گام سوم: cassandra
9.....	گام پنجم: clickhouse
9.....	مقدمه
9.....	توضیحات کد پایتون و ارتباط با Clickhouse
11.....	گام ششم: apache superset
11.....	مقدمه
11.....	نمودارها و موارد خواسته شده

گام سوم: cassandra

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.

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) );
    """)
```

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) );
    """)
```

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) );
    """)
```

همانطور که میبینید بنده از پایتون برای اجرای کوپری‌های کاساندرا استفاده کردم و این کار را با استفاده از کتابخانه cassandra انجام داده‌ام.

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;
    """)

```

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
```

در ادامه کدهای داده‌هایی را مشاهده میکنید که تعداد میگیرد و بنا بر تعداد برای ساعت، ماه و را محاسبه میکند و بر میگرداند:

```

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}-{timestamp2.month}-{timestamp2.day}', '%Y-%m-%d')
    n_year = all_post(timestamp1 , timestamp2)

    return n_hour, n_day, n_week, n_month, n_year

```

```

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}-{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

```

```

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

```

در کدها پیشبینی بزرگ شدن عدد ماه یعنی بزرگتر از 12 یا کوچکتر از 1 شده است.

گام پنجم: clickhouse

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 آنلایین 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     ]
```

بخش اول کد اضافه کردن داده به Clickhouse

```

37         'id',
38         'sendTime',
39         'sendTimePersian',
40         'hashtags',
41         'key_words',
42         'senderUsername',
43         'senderName',
44         'content'
45     )'
46     'VALUES ', data)
47
48     print(jData['id'])
49     return True
50
51 def insert_jsons(data):
52     for i in data:
53         insert(i)
54
55 import json
56 # use the below code to read the data
57 f2 = open('preprocessed_data.json', encoding="utf8")
58 jdata = json.load(f2)
59 f2.close()
60 insert_jsons(jdata)

```

بخش دوم کد اضافه کردن داده به Clickhouse

```

clickhouse-client
...
[('id','برک','id','sendTime','sendTimePersian','hashtags','key_words','senderUsername','senderName','content')
VALUES ('520f188a-7ae2-4fb6-9378-f142ed3eb56f','2021-06-27 21:47:35','1400/04/07 02:17','[...]',...)]
...
5731 rows in set. Elapsed: 1.015 sec. Processed 5.73 thousand rows, 3.17 MB (5.65 thousand rows/s., 3.13 MB/s.)

```

داده ها در محیط Clickhouse

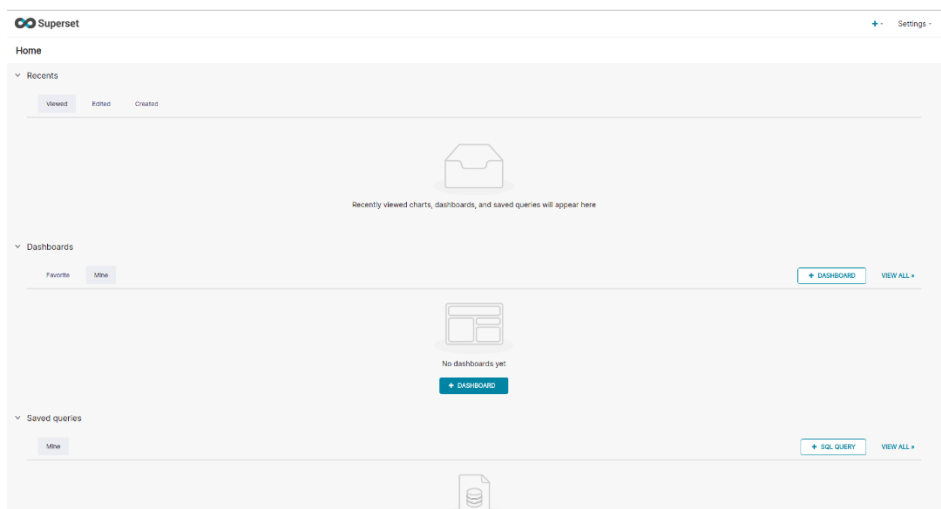
گام ششم: apache superset

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.



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

841 rows 00:00:00.57

hashtags

hashtags	COU
["u0628u0631u06a9u062a"]	1345
["u0634u0627u062fu0635_u0628u0648u0631u0633"]	512
["u0634u067u0644u06cc"]	182
["u062fu06cc"]	149
["u0634u0633u062u0627"]	115
["u0634u0646u0648u0648u0634"]	111
["u0634u067u0648u0627"]	102
["u0646u0648u0631u06cc"]	84
["u062u064fu0633u062u0631"]	71
["u0648u0627u06ccu0648u0627"]	55
["u0648u0627u06ccu0648u0642"]	48
["u067u062u0635_u0628u0648u0644"]	46
["u0634u0644u0631u0627"]	41
["u062u0648u062fu0631u0648"]	39
["u062u0632u0627u0645u06ccu0627"]	38
["u0627u06cc"]	37

▼ Data

VIEW RESULTS VIEW SAMPLES

841 rows retrieved

hashtags

hashtags	COU
["u0628u0631u06a9u062a"]	1345
["u0634u0627u062fu0635_u0628u0648u0631u0633"]	512
["u0634u067u0644u06cc"]	182
["u062fu06cc"]	149
["u0627u06cc"]	115

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
["u062u064fu0633u0627u06ccu0627_u0634u0627u062fu0633"]	10
["u0628u0648u0631u0633_u0645u0642u0627u06ccu0633u0647_u062fu0631u0635u0627u06cc"]	9
["u0627u0641u0632u0627u06ccu0634_u0633u063fu0649u0627u06ccu0647"]	8
["u0631u06ccu0627u0644_u0645u0627u0647u0647"]	7
["u0641u0631u0648u0634_u0645u0627u0647"]	7
["u0645u0628u0644u0634_u0631u06ccu0627u0644"]	6
["u0627u0644u0627u0631_u0627u0645u0627u0631u0627u062a_u0627u06ccu0631u0627u0646"]	6
["u0688u06a9u0627u0631u0646_u010"]	6
["u0637u0644u0627_u0648u0633u0628u062a_u0627u0648u0631u0647"]	6
["u062u0631u0648u0627u06cc_u0627u0627u0628u0631"]	5
["u0628u0648u0631u0633_u0648u0645u0627u062fu0647u0627_u0641u0631u0648u0634"]	5
["u0631u0648u0648u0645u0627u06ccu06cc_u0648u0627u06a9u0633u0648"]	4
["u0627u0627u06ccu0632u0647_u0627u0636u0627u0641u0647"]	4
["u0628u0648u0631u0633_u0645u0627u0647_u062u0633u0647u06ccu0644u0627u062a"]	4

▼ Data

VIEW RESULTS VIEW SAMPLES

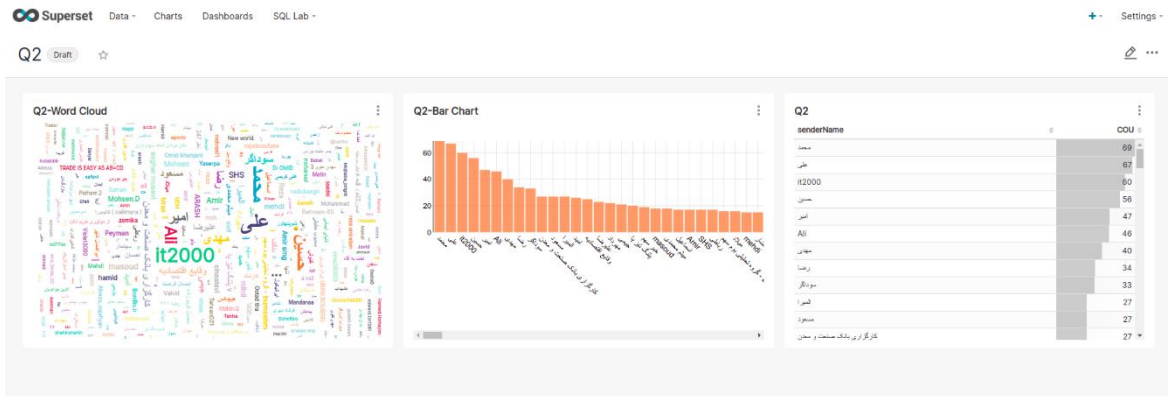
1k rows retrieved

key.words

key.words	COU
["010"]	27
["u0628u0631u06a9u062a_u0645u0634u06a9u06cc"]	24
["u062u064fu0633u0627u06ccu0627_u0634u0627u062fu0633"]	10
["u0628u0648u0631u0633_u0645u0642u0627u06ccu0633u0647_u062fu0631u0635u0627u06cc"]	9
["u0627u0641u0632u0627u06ccu0634_u0633u063fu0649u0627u06ccu0647"]	8

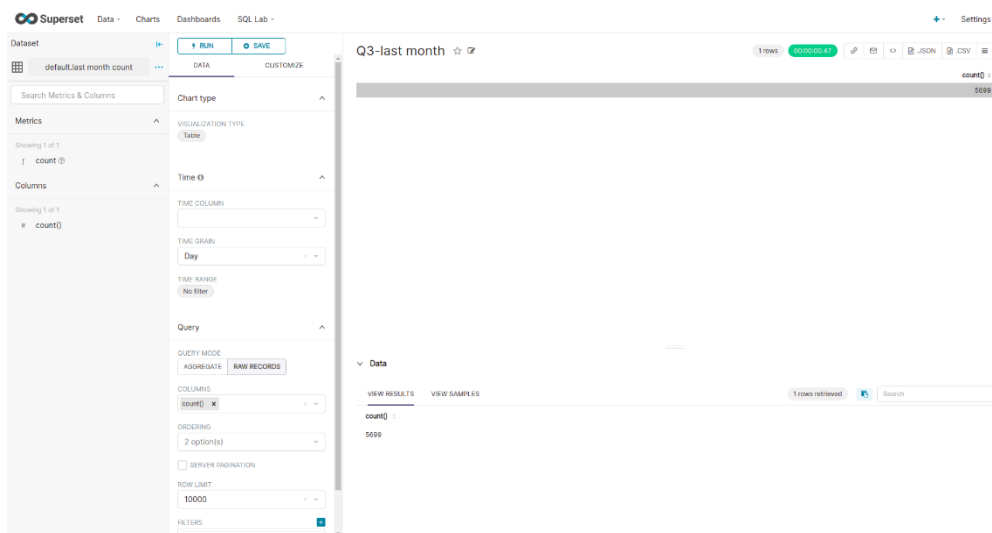
کلمات کلیدی به همراه تعداد و نمودار آنها به صورت اسکی در محیط superset

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.

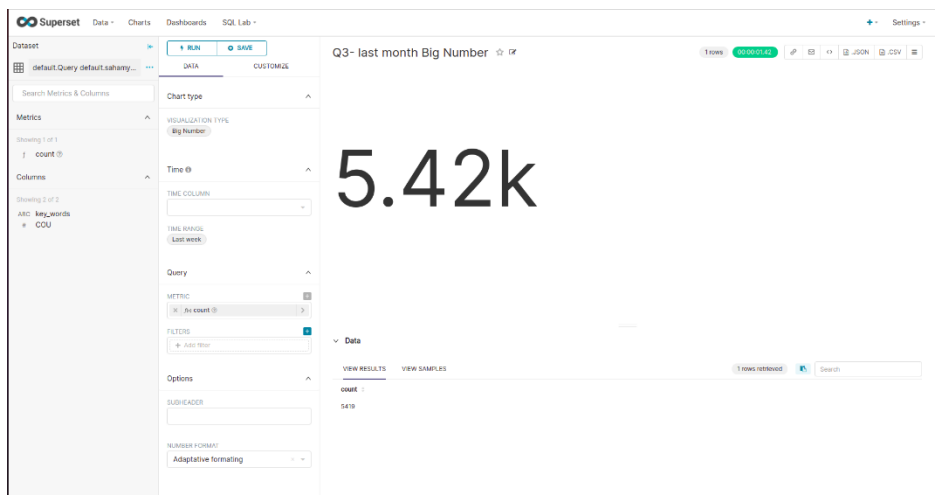


تعداد توییت های کاربران به صورت نزولی در قالب جدول نمودار و ابر کلمات در محیط superset

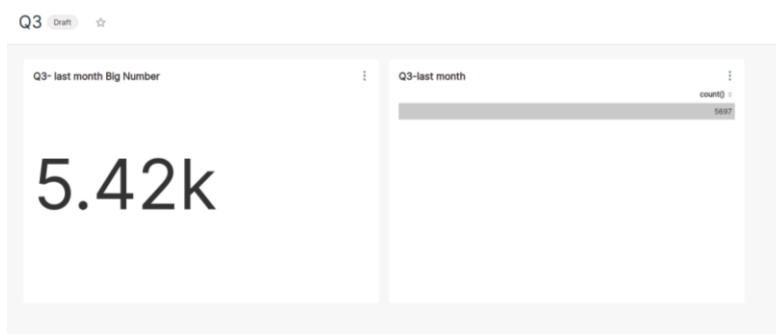
در ادامه تعداد توییت ها را در زمان ها مختلف بررسی کردیم که در تصویر 33 تعداد به دست آمده برای ماه اخیر را میبینیم و در تصاویر بعدی فرم این داده را به صورت Big Number کشیدیم.



تعداد توییت ها در ماه اخیر در محیط superset



superset Big Number تعداد توییت ها در ماه اخیر در محیط



superset Big Number و تعداد توییت ها در ماه اخیر در محیط