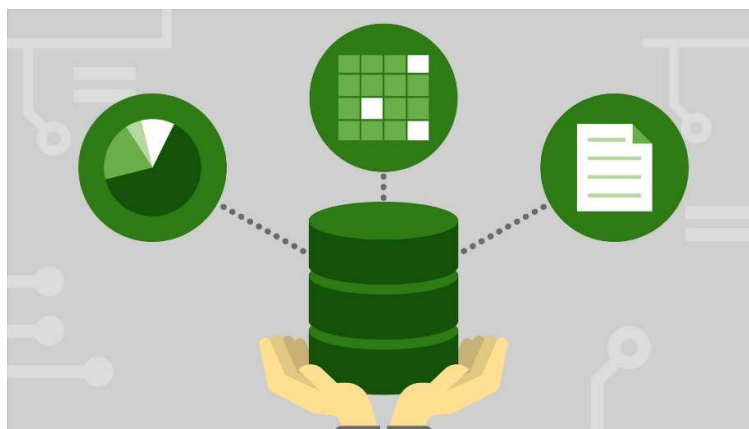


به نام خدا



دانشگاه تهران
پردیس دانشکده‌های فنی
دانشکده برق و کامپیوتر



آزمایشگاه پایگاه داده

دستورکار شماره ۷

مهلت تحویل :

۱۴۰۰/۲/۳۱

مجتبی بنائی

دستور کار شماره ۲ - کار با الاستیک سرچ

یکی از دیتابیس های معروف در حوزه جستجوی متن با سرعت بالا در ذخیره انواع داده های متنی و پاسخگویی به انواع کوئری های کاربر بر روی آنها، الاستیک سرچ است که در اکوسیستم استارتآپی ایران هم بسیار پرتعداد است.

در این دستور کار هم مشابه با دستور کارهای اخیر، هدف اصلی، آشنایی اولیه با این دیتابیس و نحوه کار با آن است که حداکثر با دوساعت صرف زمان، میتوانید به راحتی آنرا انجام دهید.

کافی است به آدرس زیر مراجعه کرده و تمامی مراحل آنرا انجام دهید :

<http://yun.ir/fi9loe>

تمام دستورات آنرا از ابتدا تا انتها در محیط کیبانا که محیط گرافیکی کار با الاستیک سرچ است انجام داده، با گرفتن اسکرین شات از خروجی آنها، گزارش خود را آپلود کنید. داده ها را طوری وارد کنید که هر کوئری حداقل دوجواب در خروجی برگرداند.

در انتهای کار، با صدا زدن API زیر در یک برنامه پایتون حداقل هزار توثیت را در الاستیک سرچ ذخیره کرده (هشتک ها که همان نمادهای بورسی هستند را جداگانه در یک لیست بریزید و سپس در الاستیک سرچ ذخیره کنید) و یک داشبورد با حداقل دو ویژوالیزیشن ایجاد کنید (مثلا ابرهشتک ها یا تعداد نمادهای پرتکرار)

<https://www.sahamyab.com/guest/twiter/list?v=0.1>

نمونه کد مورد نیاز و نحوه ایجاد یک داشبورد بر اساس داده های متنی می تواند به صورت زیر باشد :

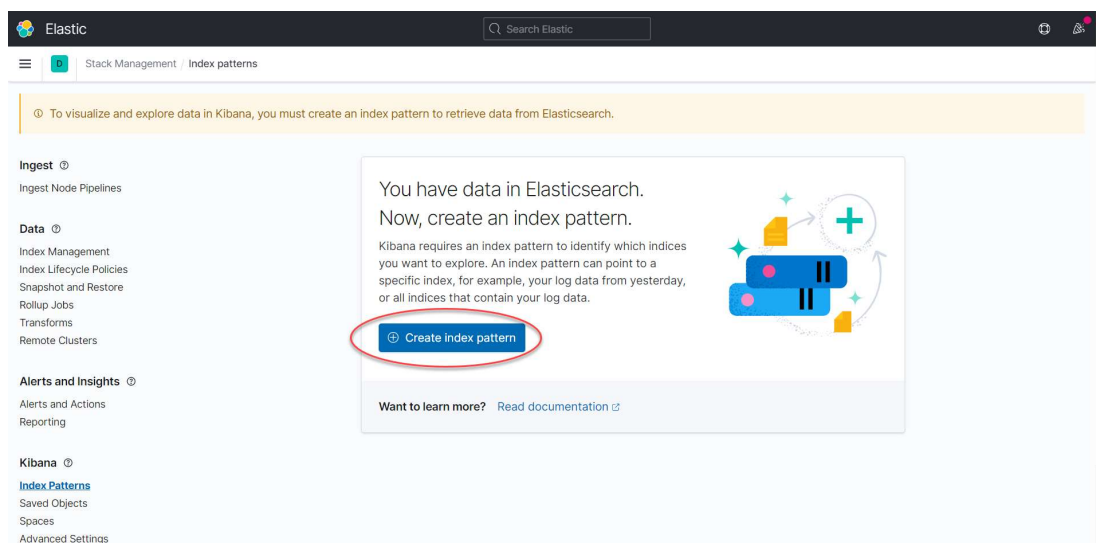
```

script.py
1 import requests
2 from elasticsearch import Elasticsearch
3 import re
4
5 url = 'https://www.sahamyab.com/guest/twitter/list?v=0.1'
6 elasticSearch = Elasticsearch([{'host': 'localhost', 'port': 9200}])
7 total = 1000
8 fetched = 0
9 seenIds = set()
10 hashtags = list()
11
12 while fetched < total:
13     response = requests.get(url=url)
14     if response.status_code != 200:
15         print('HTTP', response.status_code)
16         continue
17     data = response.json()["items"] # Check the JSON Response Content documentation below
18     for tweet in data:
19         if tweet["id"] not in seenIds:
20             try:
21                 tweet["hashtags"] = re.findall(r"#(\w+)", tweet["content"])
22                 elasticSearch.index(index="twitter", doc_type="twitter", body=tweet)
23                 seenIds.add(tweet["id"])
24                 fetched += 1
25                 print("tweet " + str(tweet["id"]) + " fetched, total: " + str(fetched))
26             except Exception as e:
27                 print(e)

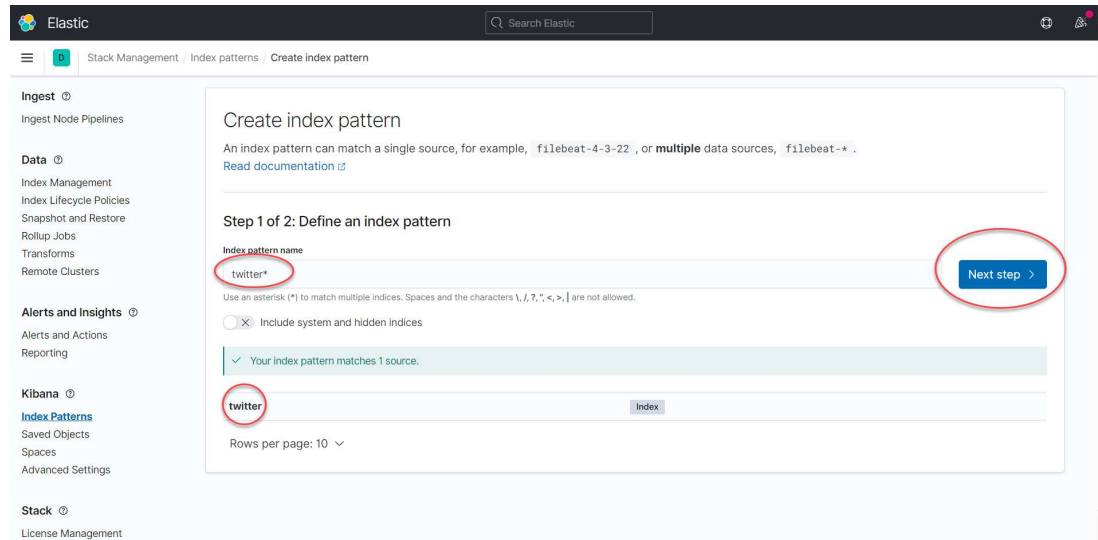
```

نحوه ایجاد داشبورد هم به صورت زیر است :

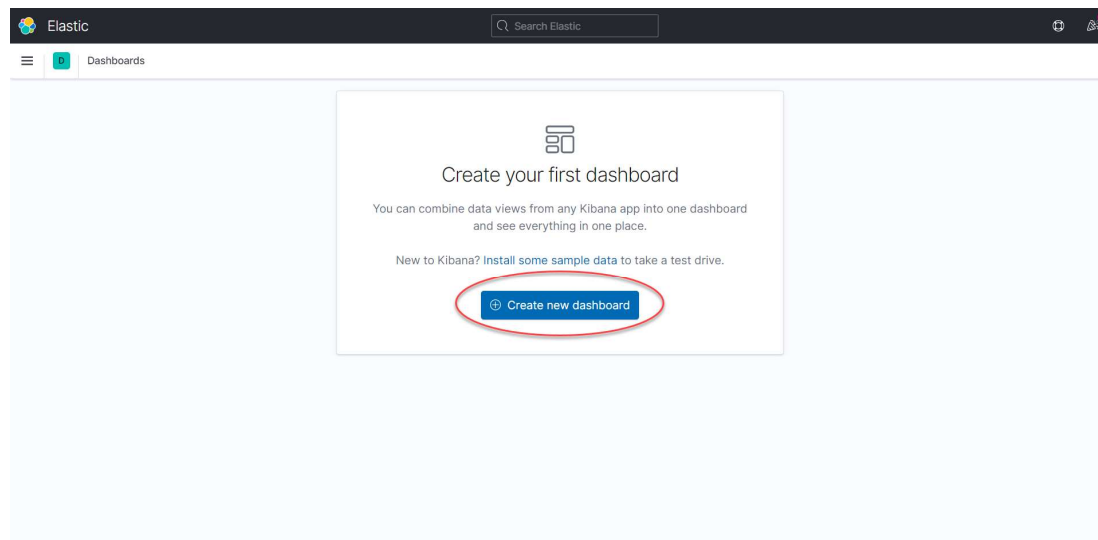
یک index pattern می سازیم.



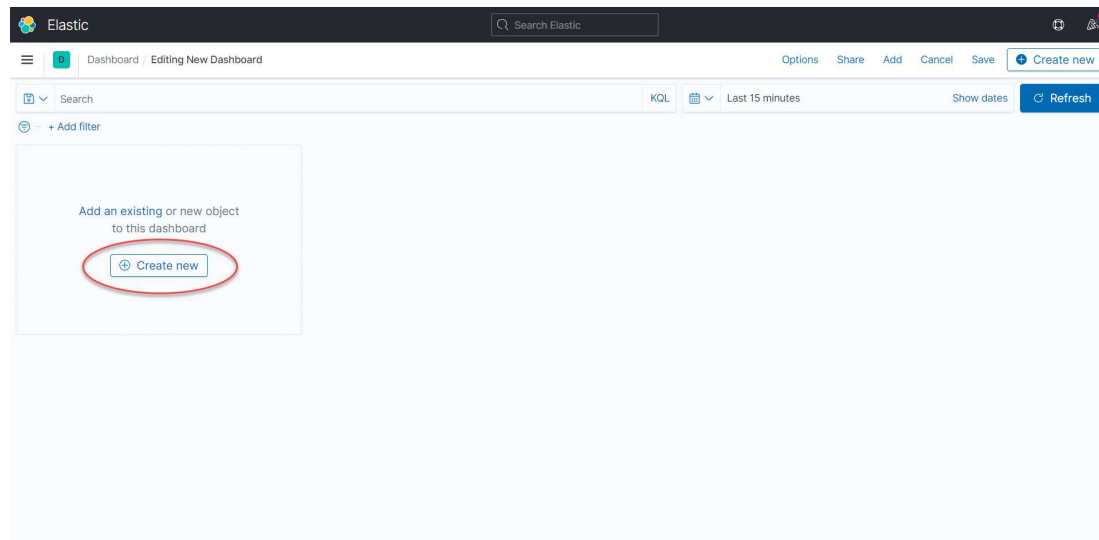
ایندکس با نام twitter را انتخاب کرده و next را می زنیم.



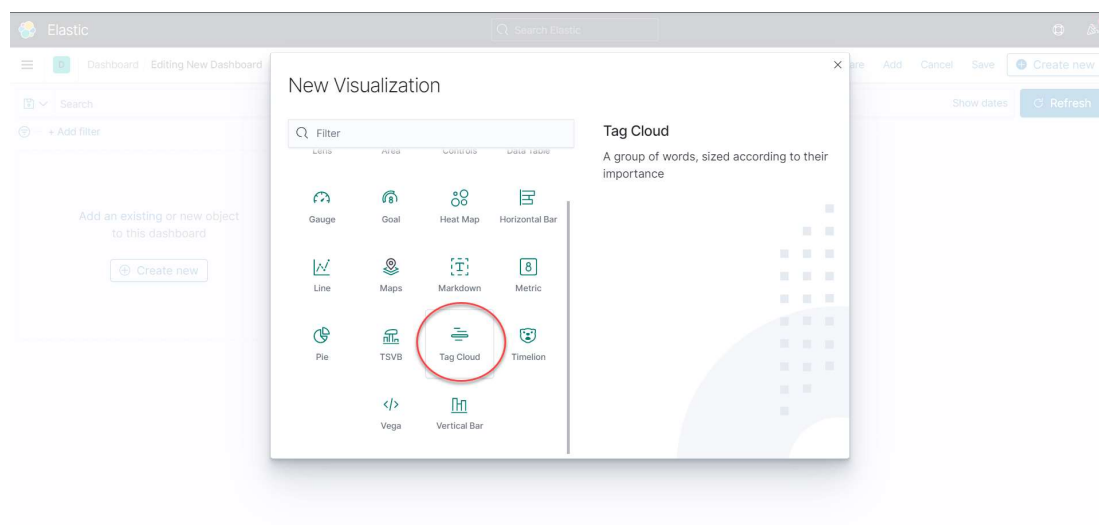
سپس یک داشبورد می سازیم.



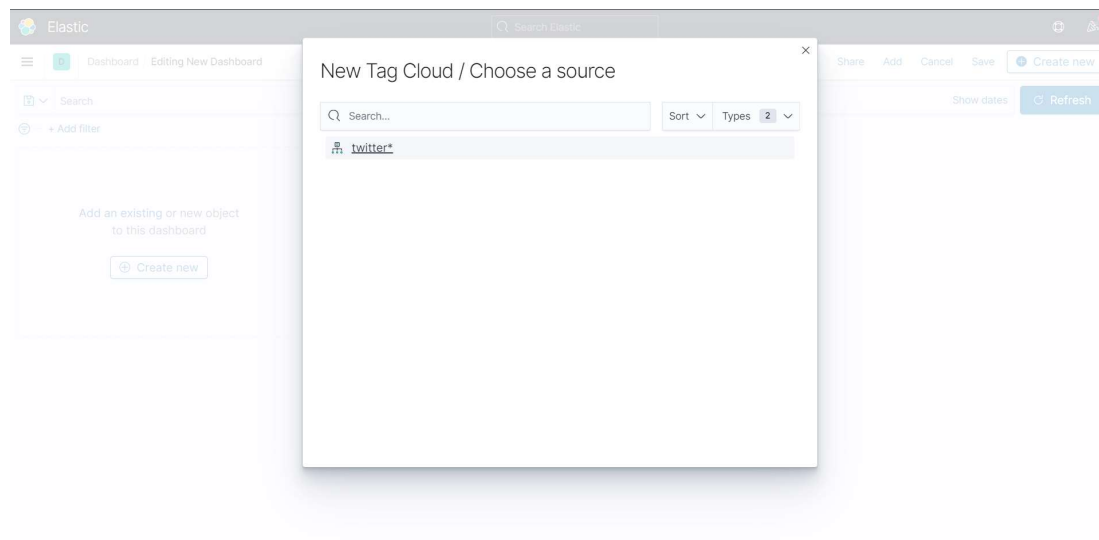
سپس برای create new را می زنیم تا twitter را اضافه کنیم.



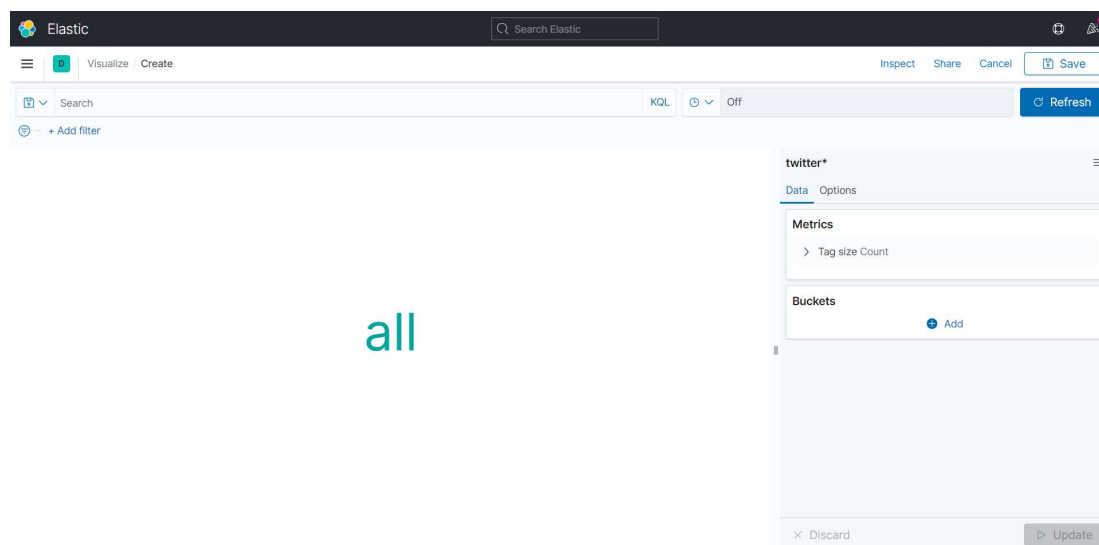
برای visualization اول tag cloud را انتخاب می کنیم که بسته به تعداد تکرار کلمه، سایز آن ست می کند.



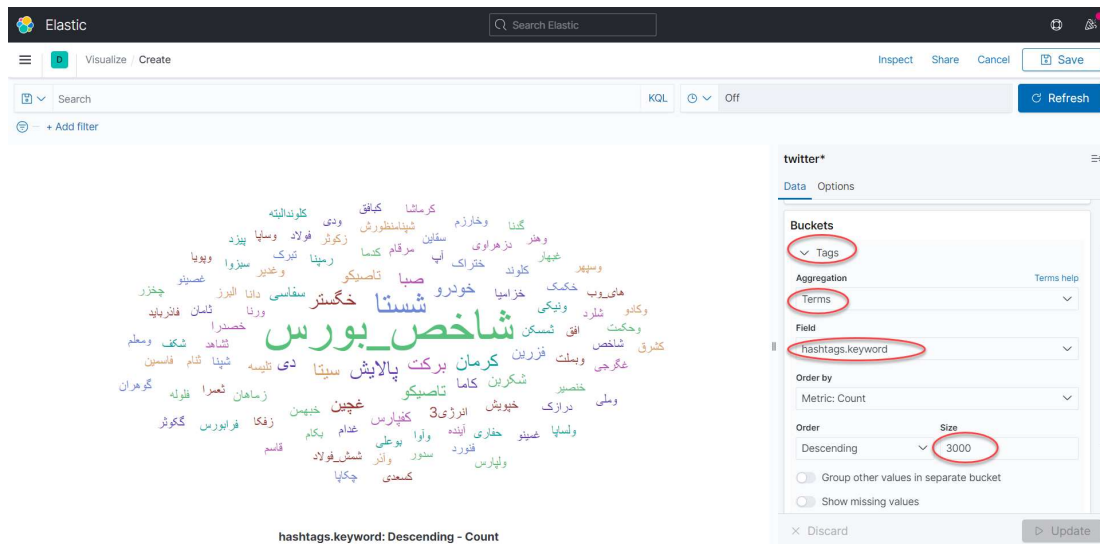
سپس twitter را انتخاب می کنیم.



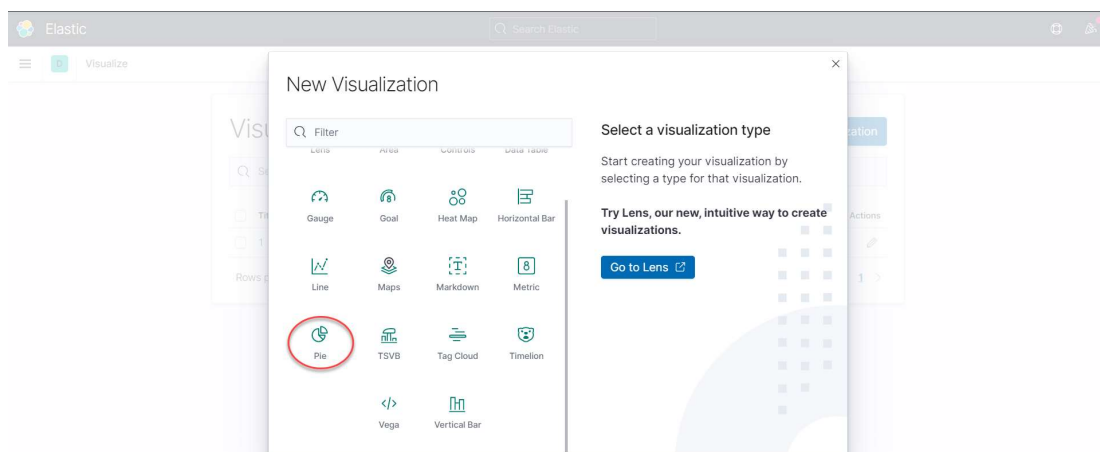
در ابتدا تنها **all** نشان داده می شود.



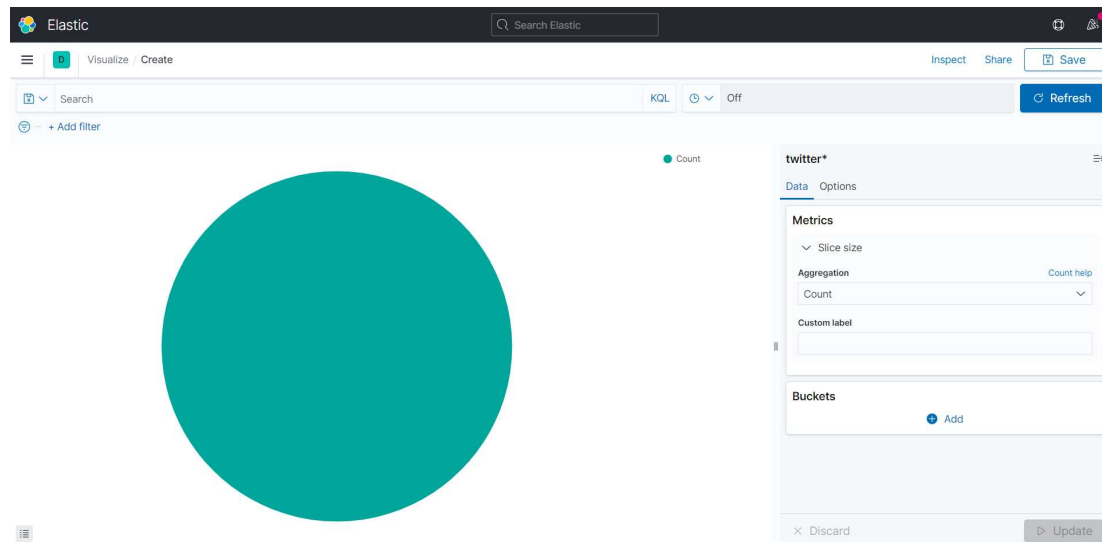
تنظیمات رو در سمت راست تغییر می دهیم تا برای **hashtag** ها این **visualization** انجام شود که خروجی به صورت زیر می شود.



برای visualization دوم، Pie را انتخاب می کنیم که سایز slice دایره بر اساس تعداد تکرار کلمه می شود.



خروجی اولیه به این صورت می شود.



در تنظیمات bucket مانند قبل hashtags رو انتخاب می کنیم و خروجی نهایی به این صورت می شود.

