

Elasticsearch як ядро пошукової СИСТЕМИ

Виконав студент прикладної математики - 3

Федусов С.В.

Науковий керівник

Доктор технічних наук, доцент Глибовець
А.М.

Вступ

- Сучасні пошукові системи мають наступні базові функції: швидкий повнотекстовий пошук, фільтрацію, сортування і ранжування документів, нечіткий пошук, швидке збереження та індексування структурованих даних.
- Необхідність отримувати результати пошукових запитів за мілісекунди, маючи десятки або сотні гігабайт інформації вимагає використання ефективних алгоритмів та структур даних. Окрім цього, сама система повинна гарантувати цілісність та відмовостійкість. Надає необхідний функціонал та задовольняє наведені вимоги пошуковий двигун Elasticsearch.

Постановка задачі

- Ознайомитися з базовим функціоналом пошукового двигуна Elasticsearch
- Дослідити алгоритми та структури даних, що використовуються для індексації
- Розглянути архітектуру пошукового двигуна
- Застосувати Elasticsearch для індексації та аналізу тестових даних

Тестові дані

- Для дослідження можливостей індексації, пошуку та аналізу даних була використана глобальна статистика поширення вірусу COVID-19

Province/State	Country/Region	Lat	Long	1/22/20	1/23/20	1/24/20	1/25/20	1/26/20
	Afghanistan	33.0	65.0	0	0	0	0	0
	Albania	41.1533	20.1683	0	0	0	0	0
	Algeria	28.0339	1.6596	0	0	0	0	0
	Andorra	42.5063	1.5218	0	0	0	0	0
	Angola	-11.2027	17.8739	0	0	0	0	0
	Antigua and Barbuda	17.0608	-61.7964	0	0	0	0	0
	Argentina	-38.4161	-63.6167	0	0	0	0	0
	Armenia	40.0691	45.0382	0	0	0	0	0
Australian Capital Territory	Australia	-35.4735	149.0124	0	0	0	0	0
New South Wales	Australia	-33.8688	151.2093	0	0	0	0	3
Northern Territory	Australia	-12.4634	130.8456	0	0	0	0	0
Queensland	Australia	-28.0167	153.4	0	0	0	0	0
South Australia	Australia	-34.9285	138.6007	0	0	0	0	0
Tasmania	Australia	-41.4545	145.9707	0	0	0	0	0
Victoria	Australia	-37.8136	144.9631	0	0	0	0	1
Western Australia	Australia	-31.9505	115.8605	0	0	0	0	0
	Austria	47.5162	14.5501	0	0	0	0	0
	Azerbaijan	40.1431	47.5769	0	0	0	0	0
	Bahamas	25.0343	-77.3963	0	0	0	0	0
	Bahrain	26.0275	50.55	0	0	0	0	0
	Bangladesh	23.685	90.3563	0	0	0	0	0
	Barbados	13.1939	-59.5432	0	0	0	0	0

This request does not have a body

Status: 200 OK Time: 80 ms Size: 343 B

Save Response



```
1 {
2   "cluster_name": "docker-cluster",
3   "status": "green",
4   "timed_out": false,
5   "number_of_nodes": 2,
6   "number_of_data_nodes": 2,
7   "active_primary_shards": 5,
8   "active_shards": 10,
9   "relocating_shards": 0,
10  "initializing_shards": 0,
11  "unassigned_shards": 0,
12  "delayed_unassigned_shards": 0,
13  "number_of_pending_tasks": 0,
14  "number_of_in_flight_fetch": 0,
15  "task_max_waiting_in_queue_millis": 0,
16  "active_shards_percent_as_number": 100.0
17 }
```

Робота з Elasticsearch

Взаємодія з кластером Elasticsearch
відбувається за допомогою http запитів.

POST http://localhost:9200/covid_19/_search

Send Save

Params Authorization Headers (10) **Body** Pre-request Script Tests Settings Cookies Code

none form-data x-www-form-urlencoded **raw** binary GraphQL **JSON** Beautify

```
1 {
2   "query": {
3     "match_all": {}
4   },
5   "from": 0,
6   "size": 10
7 }
```

POST http://localhost:9200/covid_19/_doc/

Send Save

Params Authorization Headers (10) **Body** Pre-request Script Tests Settings Cookies Code

none form-data x-www-form-urlencoded **raw** binary GraphQL **JSON** Beautify

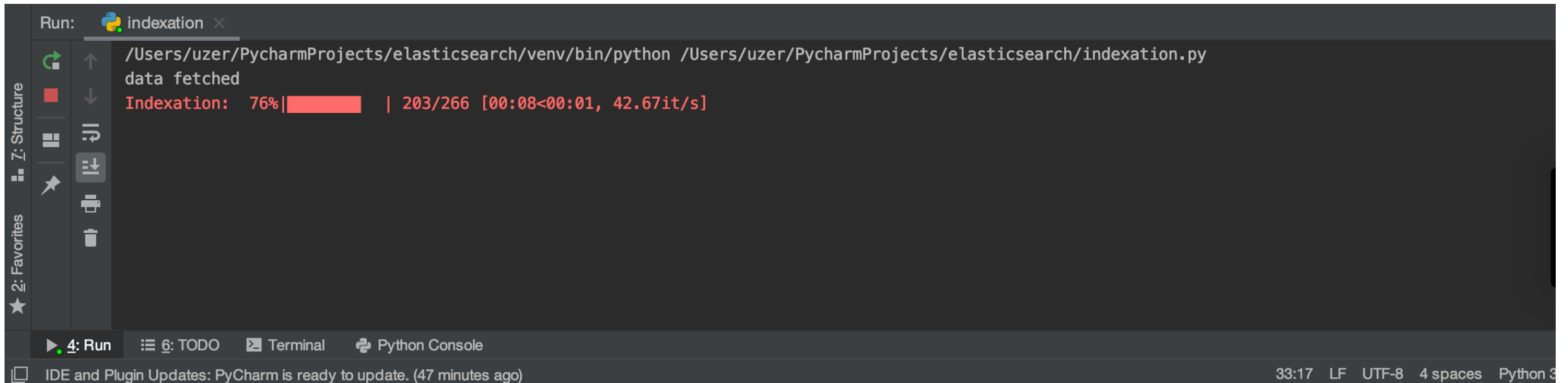
```
1 {
2   "country": "Ukraine",
3   "location": {
4     "lat": 48.3794,
5     "lon": 31.1656
6   },
7   "daily": [
8     {
9       "date": "2020-01-22",
10      "cases": 0
11    },
12    {
13      "date": "2020-01-23",
14      "cases": 0
15    },
16    {
17      "date": "2020-01-24",
18      "cases": 0
19    }
20  ]
21 }
```

Демо-застосунок для індексації

Для індексації тестових даних
було використано Python,
Pandas та Elasticsearch Python
API.



Для моніторингу процесу індексації було використано Python бібліотеку tqdm, що відображає індикатор виконання індексації, а також базову інформацію: час виконання, швидкість.



```
Run: indexation x
/Users/user/PycharmProjects/elasticsearch/venv/bin/python /Users/user/PycharmProjects/elasticsearch/indexation.py
data fetched
Indexation: 76% |███████████| 203/266 [00:08<00:01, 42.67it/s]
```

4: Run 6: TODO Terminal Python Console

IDE and Plugin Updates: PyCharm is ready to update. (47 minutes ago) 33:17 LF UTF-8 4 spaces Python 3

Розподілена архітектура

Індекс було створено з налаштуванням п'яти шард і двох реплік. Таким чином, у кластері з двох машин кожна з них містить усі шарди індексу.

Metrics

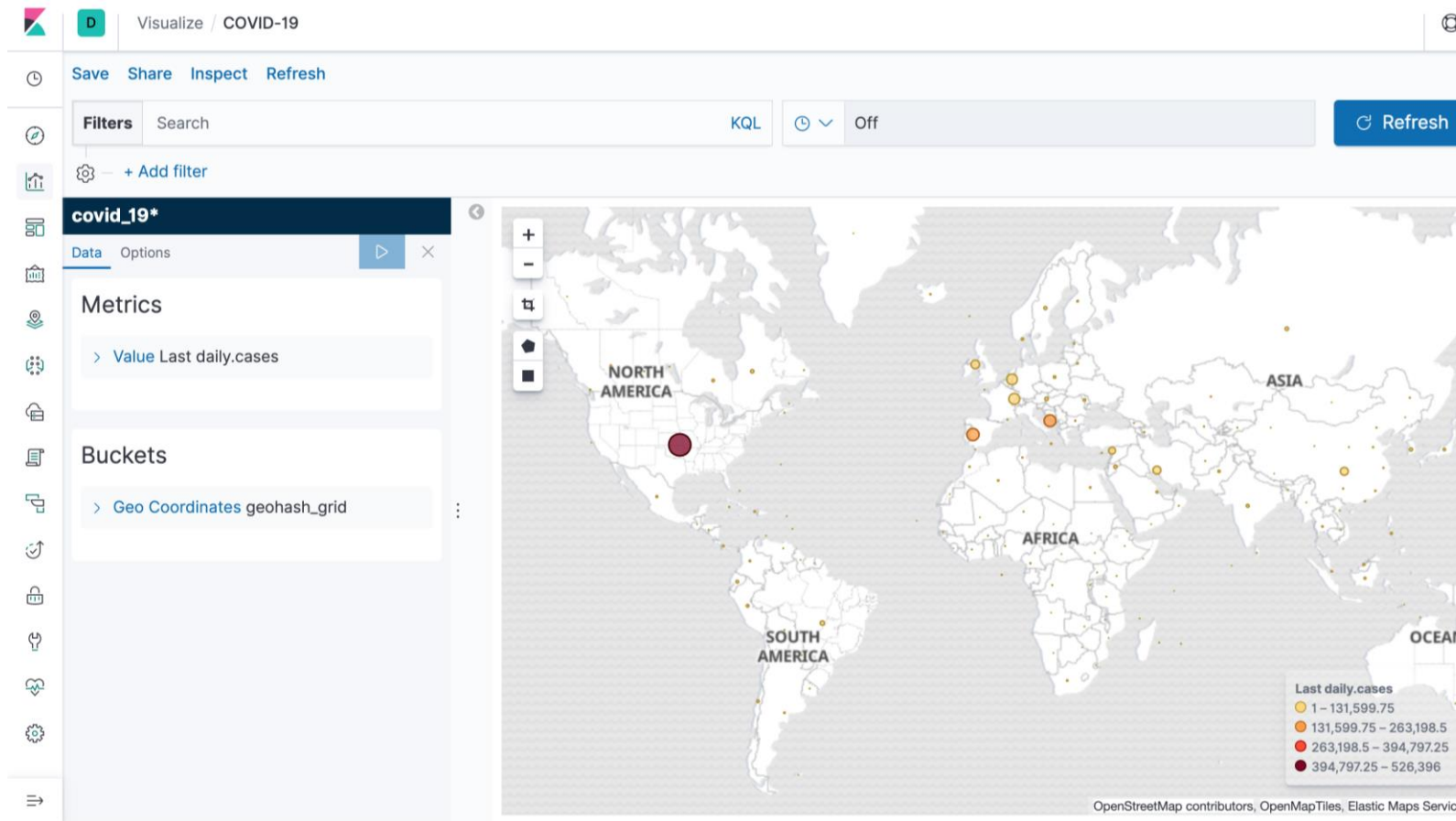
Shards

Aliases

Mappings

Administration

Shard	State	# Docs	Size	Primary	Node
0	STARTED	52	107.6kb	true	es01
0	STARTED	52	107.6kb	false	es02
1	STARTED	36	81.2kb	false	es01
1	STARTED	36	81.3kb	true	es02
2	STARTED	55	115.2kb	true	es01
2	STARTED	55	110.4kb	false	es02
3	STARTED	55	108.4kb	true	es01
3	STARTED	55	108.4kb	false	es02



Використовуючи дані розташування країн за допомогою Kibana було побудовано карту поширення вірусу. Кожна країна має відмітку у вигляді кола, що ілюструє кількість випадків захворювання

Висновки

- Схема індексу дозволяє описати будь-яку структуру, а мова запитів задовольняє будь-яку інформаційну потребу.
- Завдяки розподіленій архітектурі нівелюються фізичні обмеження однієї обчислювальної машини та забезпечується більша надійність та відмовостійкість системи.
- Elasticsearch є гарним рішенням для побудови пошукових систем

Дякую за увагу