

# Elasticsearch

## I) Introduction on Elasticsearch

I chose to work on the Elasticsearch technology. Indeed, I like distributed data systems applied to decentralized databases (NoSQL<sup>1</sup>) which is one of my specialties.

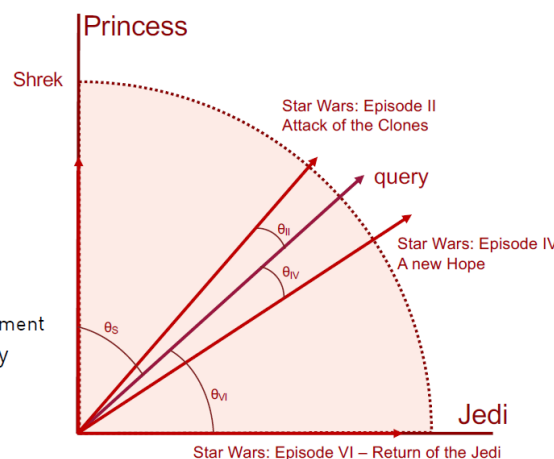
**Principle:** Elasticsearch is a powerful search engine that as its name suggests is very manageable and flexible. It has been developed on Java and is based on the DSL<sup>2</sup> language. It is however multi-language and allows a deep search and analysis of even the most complex databases.

Unlike its counterparts MongoDB, which is simpler and more versatile, or Cassandra, which is more optimized for heavier data volumes, Elasticsearch has an outstanding search quality. For this, it relies on two principles:

- **indexing:** it makes it easier to find data. Indeed, with the help of these markers that serve as shortcuts and structuring within the data, the queries are executed efficiently
- **score:** this is based on the TF<sup>3</sup> that maximizes the score and the IDF<sup>4</sup> that minimizes the score. In this way, a score is assigned to the data in order to highlight those whose score is most relevant.



- **Search Engine**
  - Similarity between
    - The query:  $q$
    - A textual document:  $d$
  - Relevance score<sup>1</sup>:  $\cos(q, d)$
- **The cosinus rely on**
  - Term Frequency (TF)
    - Normalized per key or per document
  - Inverse Document Frequency (IDF)



Practical examples of uses: Uber, Netflix, Shopify, Udemy, etc.

<sup>1</sup> **NoSQL:** databases that do not use traditional tables (rows and columns) for data storage. They organize large volumes of data using flexible techniques such as documents, charts, value pairs and columns.

<sup>2</sup> **DSL:** Domain Specific Language which is a query language.

<sup>3</sup> **TF:** term frequency.

<sup>4</sup> **IDF:** inverse document frequency.

## ELK Stack



**ELK<sup>5</sup> system:** Thus, Elasticsearch forms a trio with Logstash and Kibana.

Logstash is a data pipeline tool that is compatible with all services (S3, AWS, etc.).

Kibana is a data visualization tool.

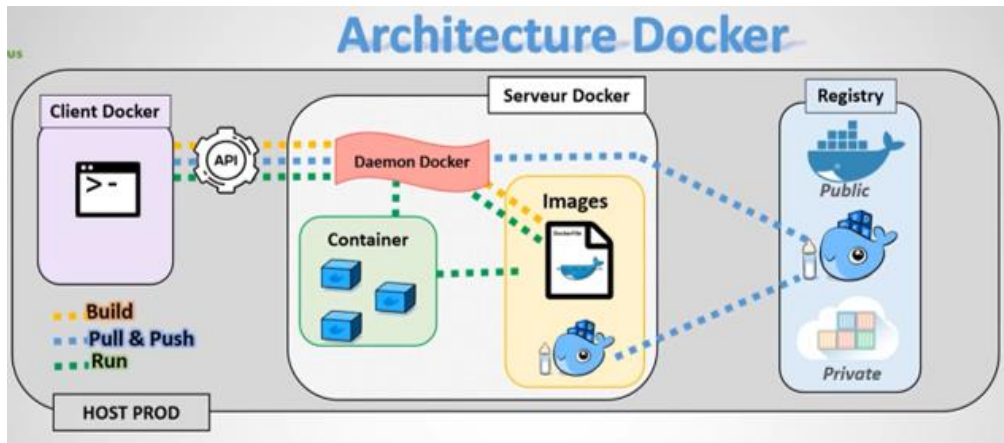
So, with this very complementary ELK system, we can import, manipulate and visualize even the most complex data efficiently.

## II) Docker

I chose to work through Docker. It's a controller based on a Linux<sup>6</sup> technology that allows you to create images and containers that work completely independently. Thus, one can create servers, connections (local) and use software in isolated environments and less impacted by the dependencies related to operating systems which allows to reduce data volumes, increase speed, portability, scalability, security and flexibility.

<sup>5</sup> **ELK:** Elasticsearch Logstash Kibana.

<sup>6</sup> **Linux:** OS exploitation system using Ubuntu technology.



In this way, I avoid installing the necessary Elasticsearch which is quite large by passing through an image that I pull from Docker. I will use an image containing Kibana to simplify the writing of queries through its interface.

```
Microsoft Windows [version 10.0.19044.2006]
(c) Microsoft Corporation. Tous droits réservés.

C:\Users\hamza>docker pull nshou/elasticsearch-kibana:latest
latest: Pulling from nshou/elasticsearch-kibana
6ec7b7d162b2: Extracting [=====>] 9.437MB/27.1MB
177617b11d13: Download complete
10273812b9e3: Download complete
ac553cdb1df6: Downloading [=====>] 33.03MB/42.4MB
1020ea78dce7: Download complete
2ea87ba2ee61: Downloading [=>] 10.74MB/486.8MB
```

After that, I configure 2 local ports (9200 for Elasticsearch and 5601 for Kibana).

NAME ↑		TAG	IMAGE ID	CREATED	SIZE
alpine/git	IN USE	latest	692618a0d74d	20 days ago	43.44 MB
docker101tutorial	IN USE	latest	f22d584f7541	18 minutes ago	28.9 MB
nshou/elasticsearch-kibana		latest	27d008902117	9 months ago	1.27 GB

Container name

elasticsearch-kibana

A random name is generated if you do not provide one.

Ports

Host port

5601

Container port

5601/tcp

Host port

9200

Container port

9200/tcp

Volumes

Host path

...

Container path

+

Environment variables

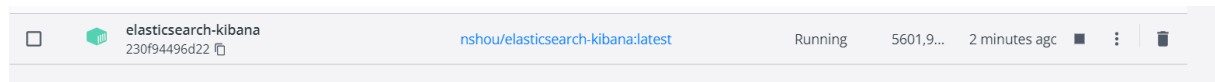
Variable

Value

+

Cancel

Run



### III) Work on the dataset

Once this is done, we can work on Elasticsearch through our localhosts. There are 3 different ways to import data:

- Logstash which I won't use.
- Curl<sup>7</sup> which allows you to use the RESTful<sup>8</sup> API.
- Kibana.

#### Curl: movies dataset

```
C:\Users\hamza\Desktop\Hamza\ESILV\2022-2023\Semestre international\UQAC\Conception et architecture des systèmes d'information\Conception - Oral I\ELASTIC\bin\elasticsearch
[2022-09-18T17:34:57.168][INFO ][o.e.n.Node               ] [BENYEMNA_HAMZA ] version[8.4.1], pid[20184], build[zip/2bd229c8e56650b42e40992322a76e7914258f0c/2022-08-26T12:11:43.232597118Z], OS[Windows 10/10.0/ma
64], JVM[Oracle Corporation/OpenJDK 64-Bit Server VM/18.0.2/18.0.2+9-g1]
[2022-09-18T17:34:57.177][INFO ][o.e.n.Node               ] [BENYEMNA_HAMZA ] JVM home [C:\Users\hamza\Desktop\Hamza\ESILV\2022-2023\Semestre international\UQAC\Conception et architecture des syst
elque\Conception - Oral I\ELASTIC\jdk], using bundled JDK [true]
[2022-09-18T17:34:57.177][INFO ][o.e.n.Node               ] [BENYEMNA_HAMZA ] JVM arguments [-Des.networkaddress.cache.ttl=60, -Des.networkaddress.cache.negative.ttl=10, -Djava.security.manager=allow, -XX:+Alway
sPreTouch, -Xss1m, -Djava.awt.headless=true, -Dfile.encoding=UTF-8, -Djna.nosys=true, -XX:-OmitStackTraceInFastThrow, -Dio.netty.noUnsafe=true, -Dio.netty.noKeySetOptimization=true, -Dio.netty.recycler.maxCapac
tyPerThread=0, -Dlog4j.shutdownHookEnabled=false, -Dlog4j2.disable.jmx=true, -Dlog4j2.formatMsgNoLookups=true, -Djava.locale.providers=SPI,COMPAT, --add-opens=java.base/java.io=ALL-UNNAMED, -XX:+UseG1GC, -Djava
io.tmpdir=C:\Users\hamza\AppData\Local\Temp\elasticsearch, -XX:+HeapDumpOnOutOfMemoryError, -XX:+ExitOnOutOfMemoryError, -XX:HeapDumpPath=data, -XX:ErrorFile=logs/hw_err_pid%p.log, -Xlog:gc*,gc+age=trace,safepo
nt:file=logs/gc.log:utctime,pid,tags:filecount=32,filesize=64m, -Xmx8135m, -Xms8135m, -XX:MaxDirectMemorySize=4205607168, -XX:G1HeapRegionSize=4m, -XX:InitiatingHeapOccupancyPercent=30, -XX:G1ReservePercent=15,
-Des.distribution.type=zip, --module-path=C:\Users\hamza\Desktop\Hamza\ESILV\2022-2023\Semestre international\UQAC\Conception et architecture des syst
modules=jdk.net, -Djdk.module.main=org.elasticsearch.server]
[2022-09-18T17:35:01.254][INFO ][c.a.c.i.j.JacksonVersion ] [BENYEMNA_HAMZA ] Package versions: jackson-annotations=2.13.2, jackson-core=2.13.2, jackson-databind=2.13.2.2, jackson-dataformat-xml=2.13.2, jackson
datatype-jsr310=2.13.2, azure-core=1.27.0, Troubleshooting version conflicts: https://aka.ms/azsdk/java/dependency/troubleshoot
[2022-09-18T17:35:04.365][INFO ][o.e.p.PluginsService     ] [BENYEMNA_HAMZA ] loaded module [ags-matrix-stats]
[2022-09-18T17:35:04.366][INFO ][o.e.p.PluginsService     ] [BENYEMNA_HAMZA ] loaded module [analysis-common]
[2022-09-18T17:35:04.366][INFO ][o.e.p.PluginsService     ] [BENYEMNA_HAMZA ] loaded module [constant-keyword]
[2022-09-18T17:35:04.367][INFO ][o.e.p.PluginsService     ] [BENYEMNA_HAMZA ] loaded module [data-streams]
[2022-09-18T17:35:04.367][INFO ][o.e.p.PluginsService     ] [BENYEMNA_HAMZA ] loaded module [frozen-indices]
[2022-09-18T17:35:04.368][INFO ][o.e.p.PluginsService     ] [BENYEMNA_HAMZA ] loaded module [ingest-attachment]
[2022-09-18T17:35:04.368][INFO ][o.e.p.PluginsService     ] [BENYEMNA_HAMZA ] loaded module [ingest-common]
[2022-09-18T17:35:04.369][INFO ][o.e.p.PluginsService     ] [BENYEMNA_HAMZA ] loaded module [ingest-geoip]
[2022-09-18T17:35:04.369][INFO ][o.e.p.PluginsService     ] [BENYEMNA_HAMZA ] loaded module [ingest-user-agent]
[2022-09-18T17:35:04.370][INFO ][o.e.p.PluginsService     ] [BENYEMNA_HAMZA ] loaded module [kibana]
[2022-09-18T17:35:04.370][INFO ][o.e.p.PluginsService     ] [BENYEMNA_HAMZA ] loaded module [lang-expression]
[2022-09-18T17:35:04.371][INFO ][o.e.p.PluginsService     ] [BENYEMNA_HAMZA ] loaded module [lang-mustache]
[2022-09-18T17:35:04.371][INFO ][o.e.p.PluginsService     ] [BENYEMNA_HAMZA ] loaded module [lang-painless]
```

We launch the server then we import the database.

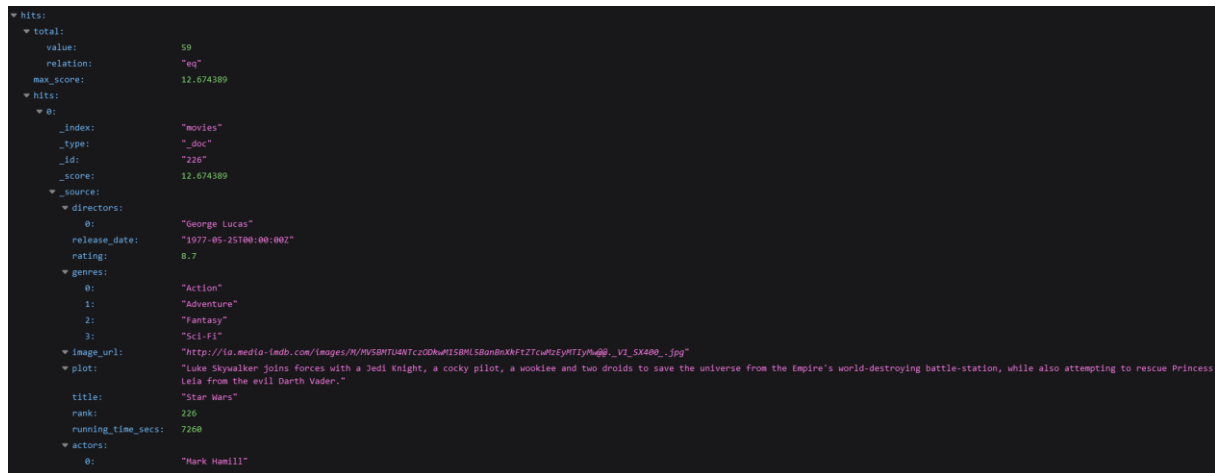
```
C:\Users\hamza\Desktop\Hamza\ESILV\2022-2023\Semestre international\UQAC\Conception et architecture des systèmes d'information\Conception - Oral I\curl -XPOT localhost:9200/_bulk -H "Content-Type: application/j
son" --data-binary @movies_elastic.json
{"took":1357,"errors":true,"items":[{"index":{"_index":"movies","_type":"_doc","_id":"1","_version":1,"result":"created","_shards":{"total":2,"successful":1,"failed":0},"_seq_no":64015,"_primary_term":2,"status":
201}},"_index":{"_index":"movies","_type":"_doc","_id":"2","_version":1,"result":"created","_shards":{"total":2,"successful":1,"failed":0},"_seq_no":64016,"_primary_term":2,"status":201}},"_index":{"_index":"mo
ovies","_type":"_doc","_id":"3","_version":1,"result":"created","_shards":{"total":2,"successful":1,"failed":0},"_seq_no":64017,"_primary_term":2,"status":201}},"_index":{"_index":"movies","_type":"_doc","_id":"4
","_version":1,"result":"created","_shards":{"total":2,"successful":1,"failed":0},"_seq_no":64018,"_primary_term":2,"status":201}},"_index":{"_index":"movies","_type":"_doc","_id":"5","_version":1,"result":"crea
ted","_shards":{"total":2,"successful":1,"failed":0},"_seq_no":64019,"_primary_term":2,"status":201}},"_index":{"_index":"movies","_type":"_doc","_id":"6","_version":1,"result":"created","_shards":{"total":2,"su
ccessful":1,"failed":0},"_seq_no":64020,"_primary_term":2,"status":201}},"_index":{"_index":"movies","_type":"_doc","_id":"7","_version":1,"result":"created","_shards":{"total":2,"successful":1,"failed":0},"_seq
_no":64021,"_primary_term":2,"status":201}},"_index":{"_index":"movies","_type":"_doc","_id":"8","_version":1,"result":"created","_shards":{"total":2,"successful":1,"failed":0},"_seq_no":64022,"_primary_term":2,"
status":201}},"_index":{"_index":"movies","_type":"_doc","_id":"9","_version":1,"result":"created","_shards":{"total":2,"successful":1,"failed":0},"_seq_no":64023,"_primary_term":2,"status":201}},"_index":{"_in
dex":"movies","_type":"_doc","_id":"10","_version":1,"result":"created","_shards":{"total":2,"successful":1,"failed":0},"_seq_no":64024,"_primary_term":2,"status":201}},"_index":{"_index":"movies","_type":"_doc
","_id":"11","_version":1,"result":"created","_shards":{"total":2,"successful":1,"failed":0},"_seq_no":64025,"_primary_term":2,"status":201}},"_index":{"_index":"movies","_type":"_doc","_id":"12","_version":1,"re
```

After that, we could make some query in the HTTP with a URL based syntax.

**Command:** [http://localhost:9200/movies/\\_search?q=title:Star+Wars](http://localhost:9200/movies/_search?q=title:Star+Wars)

<sup>7</sup> **Curl:** A HTTP request is made up of several components such as the URL to make the request to, HTTP verbs (GET, POST, etc.) and headers.

<sup>8</sup> **RESTful API:** we could totally work on a localhost browser and make some query on the dataset (<http://localhost:9200/>).



We could see the movie Star Wars with a max score of 12.67 because it is the exact name of the movie.

## Kibana: Basketwomen dataset

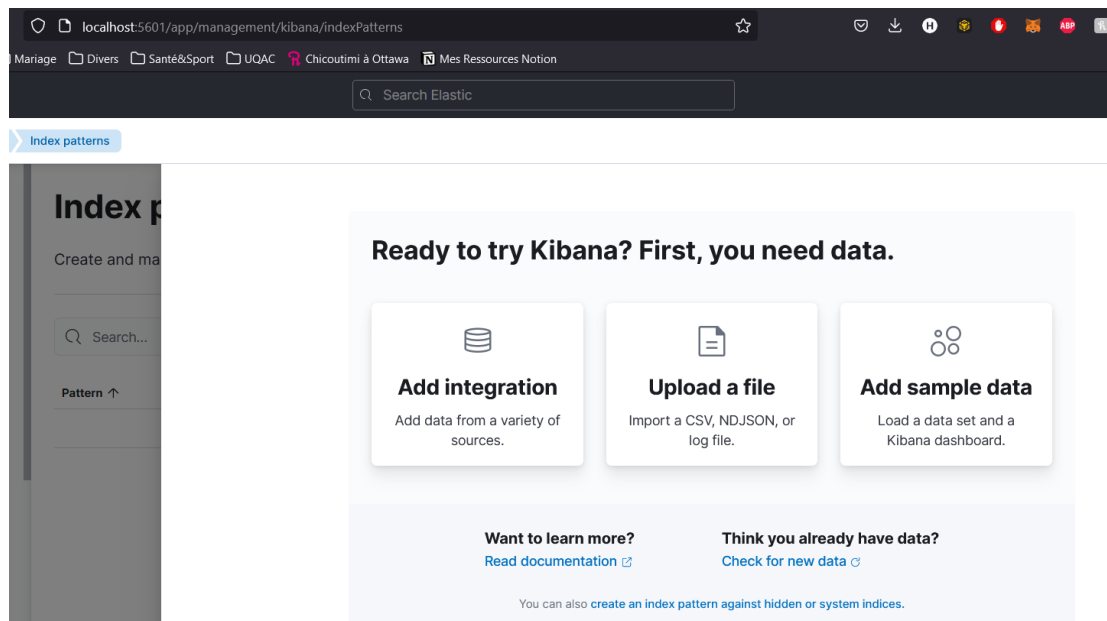
This time, I worked on Kibana which allows to write more or less complex queries with data import.

I first cleaned and corrected the format of the database required by Elasticsearch (`{"index": {"_index": "INDEXNAME" "_id": X}}`) using a Python program.

```
1 # -*- coding: utf-8 -*-
2 """
3 Created on Wed Jan 26 11:56:19 2022
4
5 @author: grego
6 """
7
8
9 MyFile= open("Basketball_women.json", 'r').read()
10 ClearData = MyFile.splitlines(True)
11 i=1
12 json_str=""
13 docs ={}
14 with open('outfile.json', 'w+') as outfile:
15     for line in ClearData:
16         #outfile.write(index_line+line)
17         outfile.write('{"index": {"_index": "Players", "_type": "Player", "_id": '+str(i)+'}}\n'+line)
18         i=i+1
19
20
21
```

Then, we import it in Kibana by going to this URL:

<http://localhost:5601/app/management/kibana/indexPatterns>

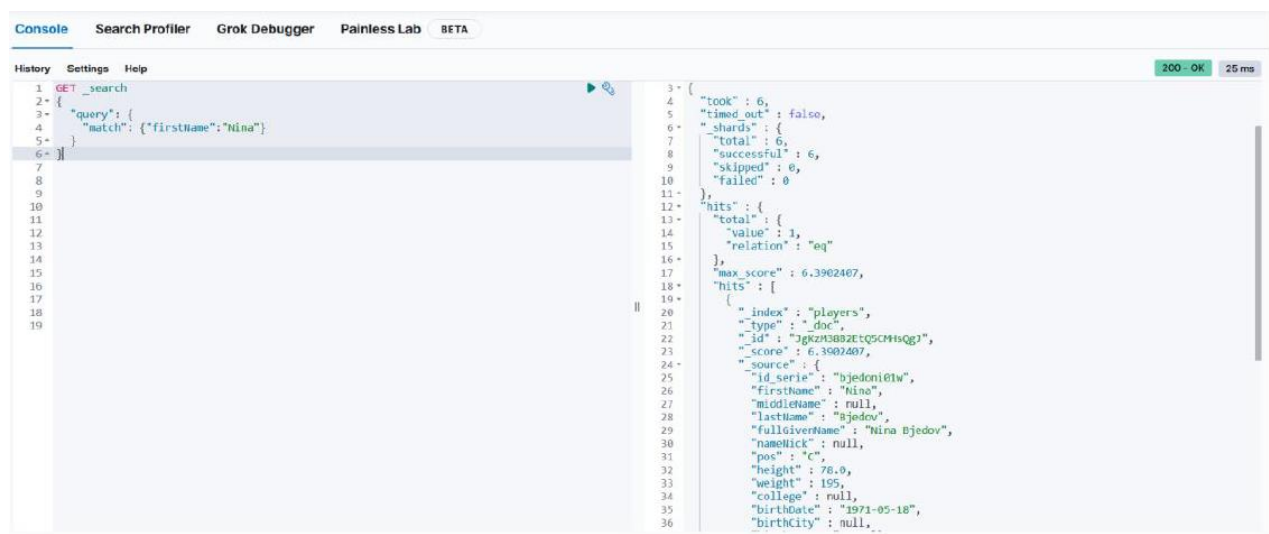


☒ Create index pattern

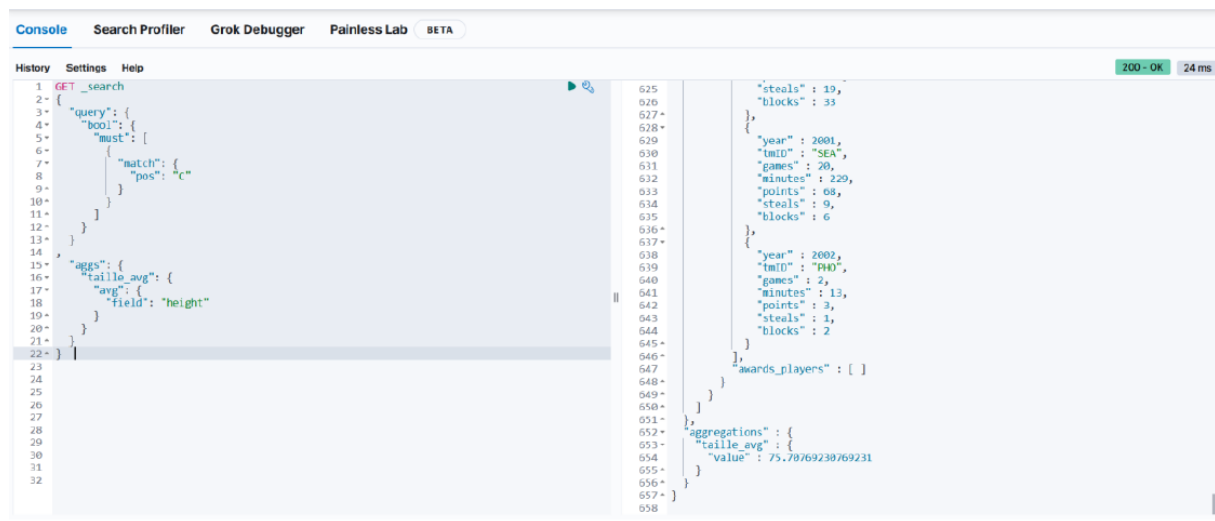
Import

The import time was quite long and difficult, but the queries are more simple to do.

**Query 1 (simple):** Basketball women named 'Nina'.



### **Query 2 (difficult with aggregation<sup>9</sup>): Average height of the basketball women playing as Center.**



```
1 GET _search
2 {
3   "query": {
4     "bool": {
5       "must": [
6         {
7           "match": {
8             "pos": "C"
9           }
10        }
11      ]
12    }
13  },
14  "aggs": {
15    "taille_avg": {
16      "avg": {
17        "field": "height"
18      }
19    }
20  }
21 }
22
```

```
625   "steals": 19,
626   "blocks": 33
627 },
628 {
629   "year": 2001,
630   "tmID": "SEA",
631   "games": 20,
632   "minutes": 229,
633   "points": 68,
634   "steals": 9,
635   "blocks": 6
636 },
637 {
638   "year": 2002,
639   "tmID": "PHO",
640   "games": 2,
641   "minutes": 13,
642   "points": 3,
643   "steals": 1,
644   "blocks": 2
645 },
646 ],
647 "awards_players": [ ]
648 },
649 },
650 },
651 },
652 },
653 },
654 },
655 },
656 },
657 },
658 }
```

```
625   "steals": 19,
626   "blocks": 33
627 },
628 {
629   "year": 2001,
630   "tmID": "SEA",
631   "games": 20,
632   "minutes": 229,
633   "points": 68,
634   "steals": 9,
635   "blocks": 6
636 },
637 {
638   "year": 2002,
639   "tmID": "PHO",
640   "games": 2,
641   "minutes": 13,
642   "points": 3,
643   "steals": 1,
644   "blocks": 2
645 },
646 ],
647 "awards_players": [ ]
648 },
649 },
650 },
651 },
652 },
653 },
654 },
655 },
656 },
657 },
658 }
```

## IV) Conclusion:

### Assets:

Indeed, Elasticsearch is a reference in the world of search engines because it is RESTful API and therefore very robust and reliable, performing with its method of indexation and scoring very effective and allows many possibilities of database management (clusters, etc.) which makes it very scalable. We thus find a technology in which we find all the strengths of distributed systems in an optimal way.

### Drawbacks:

However, it is not a good data host, especially in terms of security. The data I imported on Kibana was accessible to everyone and therefore not protected. Moreover, Elasticsearch requires data with a precise nomenclature without which the import of our data is simply impossible. Moreover, with the creation and verification of our database, the import is very long. Finally, it is a technology that takes up a lot of space. Even when importing an image via Docker where they are usually limited to a few MB, the Elasticsearch one I imported weighs more than 1Gb.

In conclusion, Elasticsearch is a great search engine tool that can be pushed more and has become a reference in the cloud and the NoSQL and cloud world.

---

<sup>9</sup> **Aggregation:** bring together distinct elements to form a homogeneous whole.

## Bibliography:

<https://www.javatpoint.com/advantages-and-disadvantages-of-elasticsearch>

<https://www.access-it.fr/actu/elasticsearch-et-ses-avantages-pour-votre-erp/>

<https://www.compose.com/articles/how-scoring-works-in-elasticsearch/>

<https://www.elao.com/blog/dev/amelioez-pertinence-resultat-elastic-search-score>

[https://www.youtube.com/watch?v=3hmA\\_K2MroY&ab\\_channel=DIGICACTUS](https://www.youtube.com/watch?v=3hmA_K2MroY&ab_channel=DIGICACTUS)

[https://www.youtube.com/watch?v=8cH0iIGlQdE&ab\\_channel=InformatiqueSansComplexe](https://www.youtube.com/watch?v=8cH0iIGlQdE&ab_channel=InformatiqueSansComplexe)

[https://www.youtube.com/watch?v=nmDN9nf1Lw&t=14s&ab\\_channel=xavki](https://www.youtube.com/watch?v=nmDN9nf1Lw&t=14s&ab_channel=xavki)

<https://chewbii.com/elasticsearch/>

Elasticsearch.pdf, Nicolas Travel, copyright ESILV.