

TME 5 - Outils pour la RI

Il existe de nombreux outils pour faire de la recherche d’information : Lucène, Indri, Terrier, Whoosh, Elastic Search (basé sur Lucène), ... Juste pour vous donner un aperçu, nous allons explorer deux d’entre eux.

Exercice 1 – Whoosh

Whoosh (<https://whoosh.readthedocs.io>) est une librairie pur-python pour la RI qui présente de nombreuses fonctionnalités.

1. Installer la librairie (<https://pypi.org/project/Whoosh/>).
2. Parcourir la documentation et explorer l’ensemble des fonctionnalités de la librairie.
3. Indexer et interroger la collection cacmShort.txt

Exercice 2 – ElasticSearch

ElasticSearch présente l’avantage de supporter des collections importantes car il est basé sur des traitements distribués.

1. Télécharger Elastic Search (<https://www.elastic.co/fr/downloads/elasticsearch> et Kibana (<https://www.elastic.co/downloads/kibana>)
2. Décompresser et lancer ElasticSearch : `./bin/elasticsearch`. documentation : <https://www.elastic.co/guide/en/elasticsearch/reference/current/starting-elasticsearch.html>
3. Décompresser et lancer Kibana : `./bin/kibana`. Documentation : <https://www.elastic.co/guide/en/kibana/6.7/start-stop.html>
4. Parcourir et exécuter certains éléments de la documentation d’ElasticSearch (principalement DocumentAPI, SearchAPI) afin de comprendre l’indexation, la recherche, ... <https://www.elastic.co/guide/en/elasticsearch/reference/current/index.html>.
Pour information, ElasticSearch est généralement lancé sur Kibana (<http://localhost:5601>) qui est une interface de visualisation et qui permet de faire des requêtes API REST. En cliquant sur "View in console" des différents exemples, vous pourrez ensuite exécuter la requête sur Kibana.
5. Les pages principales des "How-tos" :
 - Création des index <https://www.elastic.co/guide/en/elasticsearch/reference/current/indices.html>
 - Documents <https://www.elastic.co/guide/en/elasticsearch/reference/current/docs.html>
 - Recherches <https://www.elastic.co/guide/en/elasticsearch/reference/7.0/query-dsl-query-string.html>
 - Egalement : <https://www.elastic.co/guide/en/elasticsearch/reference/current/cat.html>

Il est également possible d’interfacer ElasticSearch avec Python. Très pratique pour faire des expérimentations et surtout si on souhaite parser des grandes collections de documents et les indexer. Ci-dessous, quelques éléments de code :

```
1#### Instalation
2pip3 install elasticsearch
3
4#### Usage
5####Once you started ElasticSearch, you can instantiate an ElasticSearch object :
6import elasticsearch as es
7engine = es.Elasticsearch()
8
9
```

```
10####To create an index :
11 engine.indices.create(f"{index_name}", body=settings)
12
13####Here is a complete example
14 index_name = "test-toto"
15 b = 0.5
16 k1 = 1
17 settings = {"settings": {
18     "number_of_shards": 1,
19     "index": {
20         "similarity": {
21             "default": {
22                 "type": "BM25",
23                 "b": b,
24                 "k1": k1}
25         }
26     },},
27     "mappings": {
28         "trec": {
29             "properties": {
30                 "title": {
31                     "type": "text",
32                     "index": "false"
33                 },
34                 "text": {
35                     "type": "text",
36                     "analyzer": "english"
37                 }
38             }
39         }
40     }
41 }
42 engine.indices.create(f"{index_name}", body=settings)
43
44
45####Using bulk to index multiple documents:
46 import elasticsearch.helpers as helpers
47 def bulking(texts, titles):
48     bulk_data = []
49     for i, (txt, title) in enumerate(zip(texts, titles)):
50         data_dict = {
51             '_index': 'test-toto',
52             '_type': 'trec',
53             '_id': str(i),
54             '_source': {
55                 "text": txt,
56                 "title": title
57             }
58         }
59         bulk_data.append(data_dict)
60     return bulk_data
61
62 texts = ['toto_is_now_gone',
63          'toto_has_left_the_place',
64          'toto_went_to_his_parent's_place',
65          'toto's_place_is_quite_far_away']
66
67 titles = ["Toto's_life",
68           "Toto_leaving",
```

```
69         "Toto's_trip",
70         "Toto_is_far"]
71
72 helpers.bulk(engine, bulking(texts, titles))
73
74###Search example :
75 engine.search("test-bm25", "trec", {"query": {"match": {"text": "toto_place"}}})
```

Une autre façon de faire est d'utiliser le shell : <https://openclassrooms.com/fr/courses/4462426-maitrisez-les-4474691-etudiez-le-fonctionnement-d-elasticsearch>

A vous de jouer!!!