# CE706 - Information Retrieval SU 2022

## Assigment 1

Student ID: 2112102

# Instructions for running your system

I have Mentioned how I have performed various tasks of information retrival by using online platforms elastic search and kibana, which are most commonly used softwares in industry. I have followed my lab work documents and had researched some to get better idea of these platforms. I have downloaded the dataset from the given site https://research.signal-ai.com/newsir16/signal-dataset.html. It is a signal media in which there are one million news articles gathered from different sources. After obtaining data I read that jsonl file in my code and then started performing the tasks. for that, I have splitted data into equal parts where each set had 1000 phrases only(for easy and fast analysing). After that I have started performing the tasks as asked eg. Data importing, Indexing, Tokenization, Stemming, Searching etc.

## 1. Indexing

With the help of indexing we can separate and retrieve the data from our dataset. Another concept is about replicas and shards. It is a mapping and it defines multiple types. Here I have got different indexing results using GET and PUT method.

### 1.1 By using GET method:

```
1  GET karan

 1 {
 2     "karan" : {
 3         "aliases" : { },
 4         "mappings" : {
 5             "_doc" : {
 6                 "properties" : {
 7                     "@timestamp" : {
 8                         "type" : "date"
 9                     },
10                     "content" : {
11                         "type" : "text"
12                     },
13                     "id" : {
14                         "type" : "keyword"
15                     },
16                     "media-type" : {
17                         "type" : "keyword"
18                     },
19                     "published" : {
20                         "type" : "date"
21                     },
22                     "source" : {
23                         "type" : "text"
24                     },
25                     "title" : {
26                         "type" : "text"
27                     }
28                 }
29             }
30         },
31         "settings" : {
32             "index" : {
33                 "creation_date" : "1655421473352",
34                 "number_of_shards" : "5",
35                 "number_of_replicas" : "1",
36                 "uuid" : "1SO35tlDRcydoDM1waIwJA",
37                 "version" : {
38                     "created" : "6050199"
39                 },
40                 "provided_name" : "karan"
41             }
42         }
43     }
44 }
45
```

## 1.2 By using PUT:

```
1   GET karan
2   PUT /karan1
3
4 ▾ {
5
6
7
8 ▾     "mappings":{
9
10
11
12 ▾         "properties":{
13
14
15
16 ▾             "content":{
17
18
19
20 ▾                 "source" :{
21
22
23
24                     "type" : "text" ,
25
26
27
28                     "my_signal_news":"simple"
29
30 ▴ }
31
32 ▴ }
33
34 ▴ }
35
36 ▴ }
37
38 ▴ }
39
```

```
1   #! Deprecation: the default number of shards will change from [5]
        this on the create index request or with an index template
2 ▾ {
3     "acknowledged" : true,
4     "shards_acknowledged" : true,
5     "index" : "karan1"
6 ▴ }
7
```

## 1.3 By using POST method

```
17 ▴ }
18 ▴ }
19 ▴ }
20 ▴ }
21 ▴ }
22
23   POST /karan1/_doc
24 ▾ {
25     "source":"Credit Cards"
26 ▴ }
27
28   PUT karan3
29 ▾ {
30
31 ▾ "settings": {
32 ▾ "analysis": {
33 ▾ "analyzer": {
34 ▾ "my_signal_news": {
35   "type": "standard",
36   "stopwords": "_decision_"
37 ▴ }
```

```
1 ▾ {
2     "_index" : "karan1",
3     "_type" : "_doc",
4     "_id" : "2Dwob4EBBVl0Psu5BW_H",
5     "_version" : 1,
6     "result" : "created",
7 ▾   "_shards" : {
8       "total" : 2,
9       "successful" : 1,
10      "failed" : 0
11 ▴  },
12    "_seq_no" : 0,
13    "_primary_term" : 1
14 ▴ }
15
```

# 2. Tokenization and Normalisation

The tokenizer breaks the data in individual tokens and gives output of stream of tokens. We use tokenization to protect the sensitive datawhen we are making preservation of utilities. Tokenization process is different than the encryption process in which all the sensitive data is been modified & stored.

## 2.1 Uppercase:

Here I have used uppercase filter to convert all the lowercase words in the word to the uppercase words.

## 2.2 Lowercase:

Here I have used lowercase filter to cover all the uppercase words to the lowercase words.

```
 1  +
 2  PUT /karan1
 3
 4  {
 5
 6      "mappings":{
 7
 8          "properties":{
 9
10              "content":{
11
12                  "source" :{
13
14                      "type" : "text" ,
15
16                          "my_signal_news":"simple"
17  }
18  }
19  }
20  }
21  }
22
23  PUT karan3
24  {
25
26  "settings": {
27  "analysis": {
28  "analyzer": {
29  "my_signal_news": {
30  "type": "standard",
31  "stopwords": "_decision_"
32  }
33  }
34  }
35  }
36  }
37
38
39  GET karan1/_analyze           ▶  🔧
40  {
41      "tokenizer" : "standard",
42      "filter" : ["uppercase"],
43      "text" : "TEACHERS across the nation now have access to
           information designed to help them identify whether a
           student has become radicalised"
44  }
45
46
47  PUT karan/_doc/1
```

```
 1  {
 2      "tokens" : [
 3          {
 4              "token" : "TEACHERS",
 5              "start_offset" : 0,
 6              "end_offset" : 8,
 7              "type" : "<ALPHANUM>",
 8              "position" : 0
 9          },
10          {
11              "token" : "ACROSS",
12              "start_offset" : 9,
13              "end_offset" : 15,
14              "type" : "<ALPHANUM>",
15              "position" : 1
16          },
17          {
18              "token" : "THE",
19              "start_offset" : 16,
20              "end_offset" : 19,
21              "type" : "<ALPHANUM>",
22              "position" : 2
23          },
24          {
25              "token" : "NATION",
26              "start_offset" : 20,
27              "end_offset" : 26,
28              "type" : "<ALPHANUM>",
29              "position" : 3
30          },
31          {
32              "token" : "NOW",
33              "start_offset" : 27,
34              "end_offset" : 30,
35              "type" : "<ALPHANUM>",
36              "position" : 4
37          },
38          {
39              "token" : "HAVE",
40              "start_offset" : 31,
41              "end_offset" : 35,
42              "type" : "<ALPHANUM>",
43              "position" : 5
44          },
45          {
46              "token" : "ACCESS",
47              "start_offset" : 36,
48              "end_offset" : 42,
49              "type" : "<ALPHANUM>",
```

# 3. Selecting Keywords

Keywords are the words which are frequently used by the user while searching in the search engine. Here I have used stop filter, It is used to eliminate stop words. Here analyzer makes the token of the given input to give better results.

```
 1  GET karan/_analyze           ▶  🔧
 2  {
 3      "tokenizer" : "standard",
 4      "filter": ["stop"]
 5      , "text": ["Additional security forces have been called in to
           control the situation in the Shamshabad area here"]
 6  }
```

```
 1  {
 2      "tokens" : [
 3          {
 4              "token" : "Additional",
 5              "start_offset" : 0,
 6              "end_offset" : 10,
 7              "type" : "<ALPHANUM>",
 8              "position" : 0
 9          },
10          {
11              "token" : "security",
12              "start_offset" : 11,
13              "end_offset" : 19,
14              "type" : "<ALPHANUM>",
15              "position" : 1
16          },
17          {
18              "token" : "forces",
19              "start_offset" : 20,
20              "end_offset" : 26,
21              "type" : "<ALPHANUM>",
22              "position" : 2
23          },
24          {
25              "token" : "have",
26              "start_offset" : 27,
27              "end_offset" : 31,
28              "type" : "<ALPHANUM>",
29              "position" : 3
30          },
31          {
32              "token" : "been",
33              "start_offset" : 32,
34              "end_offset" : 36,
35              "type" : "<ALPHANUM>",
36              "position" : 4
37          },
38          {
39              "token" : "called",
40              "start_offset" : 37,
41              "end_offset" : 43,
42              "type" : "<ALPHANUM>",
43              "position" : 5
44          },
45          {
46              "token" : "control",
47              "start_offset" : 50,
48              "end_offset" : 57,
49              "type" : "<ALPHANUM>",
```

## 4. Stemming or Morphological Analysis

Stemmer filter reduces the word and convert it to its root form. It takes care if the word variants matches during the search. It is a normalisation technique used in Neural Language Processing. Stemming helps to reduce the dimancity of data.

```
1   GET karan/_analyze
2   {
3       "tokenizer" : "standard",
4       "filter": ["stemmer"]
5       , "text": ["credit card"]
6   }
```

```
1   {
2       "tokens" : [
3           {
4               "token" : "credit",
5               "start_offset" : 0,
6               "end_offset" : 6,
7               "type" : "<ALPHANUM>",
8               "position" : 0
9           },
10          {
11              "token" : "card",
12              "start_offset" : 7,
13              "end_offset" : 11,
14              "type" : "<ALPHANUM>",
15              "position" : 1
16          }
17      ]
18  }
19
```

## 5. Searching

This search API helps to seach the item in elastic search. Here I have fired two quries and it gives result according to how many tokens are made of those words. It can also gives data of past times like which things or network was ran or searched by the user.

```
18  }
19  }
20  }
21  }
22
23  PUT karan3
24  {
25
26  "settings": {
27  "analysis": {
28  "analyzer": {
29  "my_signal_news": {
30  "type": "standard",
31  "stopwords": "_decision_"
32  }
33  }
34  }
35  }
36  }
37
38
39  GET karan1/_analyze
40  {
41      "tokenizer" : "standard",
42      "filter" : ["uppercase"],
43      "text" : "TEACHERS across the nation now have access to
            information designed to help them identify whether a
            student has become radicalised"
44  }
45
46
47  PUT karan3/_doc/1
48  {
49      "body": "Jamie Heaslip and Les Kiss spooke to the press
            following the Captain's Run in Wembley today."
50  }
51
52
53  GET _search
54  {
55      "query": {
56          "query_string": {
57              "query": "(chairperson) OR (Tuesday)",
58              "default_field": "content"
59          }
60      }
61  }
62
63
```

```
1   {
2       "took" : 268,
3       "timed_out" : false,
4       "_shards" : {
5           "total" : 46,
6           "successful" : 46,
7           "skipped" : 0,
8           "failed" : 0
9       },
10      "hits" : {
11          "total" : 1719,
12          "max_score" : 12.889124,
13          "hits" : [
14              {
15                  "_index" : "karan",
16                  "_type" : "_doc",
17                  "_id" : "sDzPboEBBVl0Psu5HkZZ",
18                  "_score" : 12.889124,
19                  "_source" : {
20                      "@timestamp" : "2015-09-22T13:50:01.000Z",
21                      "id" : "5db72fdd-8af9-469a-a50a-92c0be9afeb1",
22                      "source" : "New Delhi News.Net",
23                      "published" : "2015-09-22T13:50:01Z",
24                      "title" : "Malik demands probe over mysterious killings in Kashmir",
25                      "media-type" : "News",
26                      "content" : """
27  Jammu and Kashmir Liberation Front (JKLF) chairman Yasin Malik on Tuesday demanded an International probe over mysterious killings in Kashmir.
28
29  "People of Kashmir want to know who is behind these mysterious killings. We want an international probe over killings of innocent youth," he said.
30
31  Malik also condemned the Vishwa Hindu Parishad (VHP) enforcement on the economic blockade of Kashmir.
32
33  "We want to tell Vishwa Hindu Parishad (VHP) that we are not scared of their hollow threats," he said.
34
35  According to reports, the Jammu and Kashmir Liberation Front leader was detained with nearly two dozen activists and lodged at the Kothi Bagh police
    station today. Parveena Ahanger, the chairperson of the Association of Parents of Disappeared Persons, was also among those detained by the police.
    (ANI)
36  """
37              }
38          },
39          {
40              "_index" : "karan",
41              "_type" : "_doc",
42              "_id" : "-zzPboEBBVl0Psu5HjdW",
43              "_score" : 12.58832,
44              "_source" : {
45                  "@timestamp" : "2015-09-03T02:51:46.000Z",
46                  "id" : "4586ba67-4e9c-44aa-972a-f01b4400862a",
47                  "source" : "New Zimbabwe.com",
```