

University of Essex

DEPARTMENT OF MATHEMATICAL SCIENCES

Assignment 2: Elasticsearch & Evaluation

Submitted as part of the requirements for:

CE706-7-SP-CO: INFORMATION RETRIEVAL

Name: DIVYA ARORA

Registration Number: 1901423

Name: ARITRA GANGULY

Registration Number: 1906467

Supervisor: ALBA GARCIA SECO DE HERRERA

Date of submission (9TH APRIL 2020)

Introduction

This task is majorly focusing on Elasticsearch. It is an open source search engine known for text search and analytics. It is capable of handling complex search features and big data by storing, searching and analyzing rapidly. With the fact that it provides a distributed system and has compatibility with JSON, it is used to carry out the process of indexing [1].

We have to proceed with the task with the consideration that we are new to the organisation where they are in urgent need of a search engine that helps the employees to search the document collection [2].

1 Main Body

Our task is to convert documents into a structured index by using information retrieval models and evaluate them. We have to follow certain steps to fulfill our requirements i.e. indexing, searching, building a text collection, evaluation, engineering a complete system [2].

1.1 *Software Requirements*

There is a need to use three software for transforming the documents.

- Elasticsearch
 - Download Elasticsearch software.
 - Extract Zip file from download folder.
 - Right click on Elasticsearch batch file and run as administrator.
 - Command prompt window open and continue running.
 - To check whether it is installed properly. We need to open any web browser and use this URL: <http://localhost:9200/>

Elasticsearch is used for creating indexes for the documents.

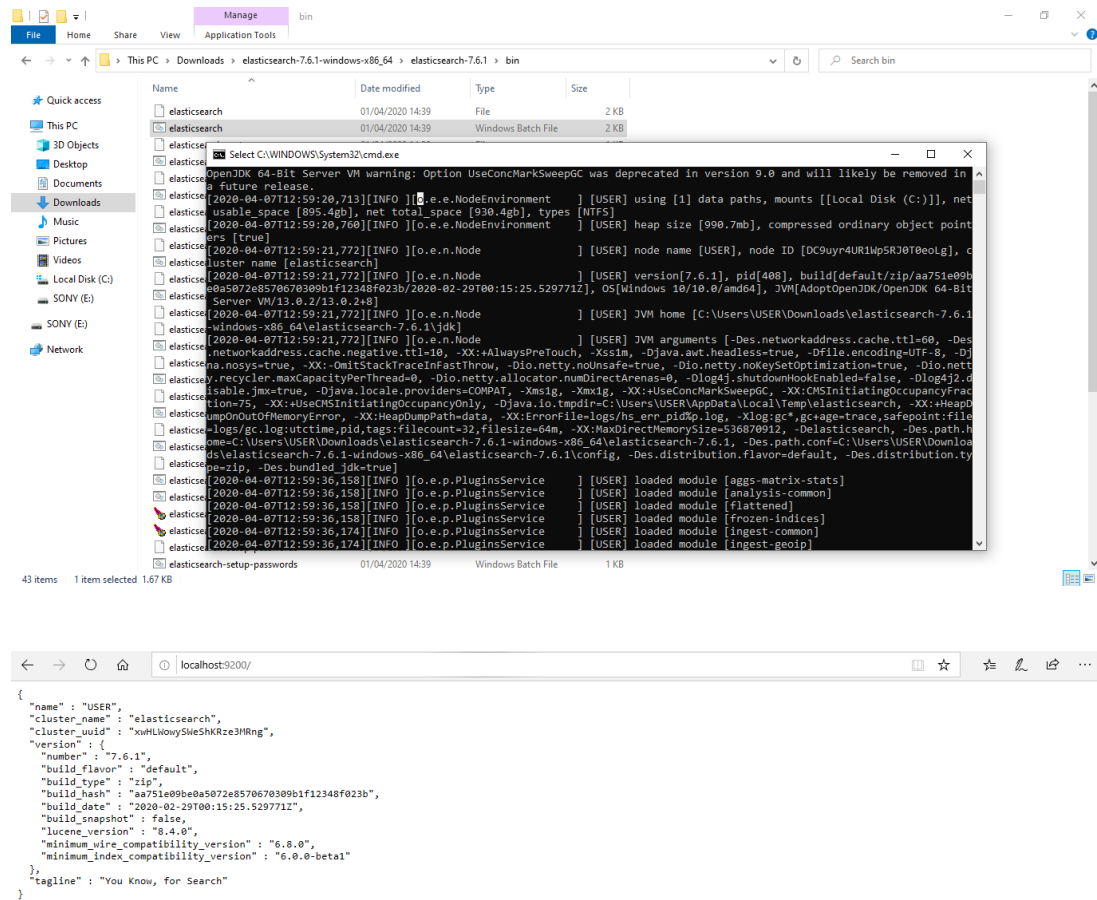


Figure 1: Running Elasticsearch

- Kibana
 - Download Kibana software.
 - Extracting Zip folder.
 - Right click on Kibana batch file and run as administrator.
 - Command prompt open showing the execution of software.

- To check whether Kibana is installed, open web browser and access URL: <http://localhost:5601/>

Kibana is used to search for the index created using Elasticsearch.

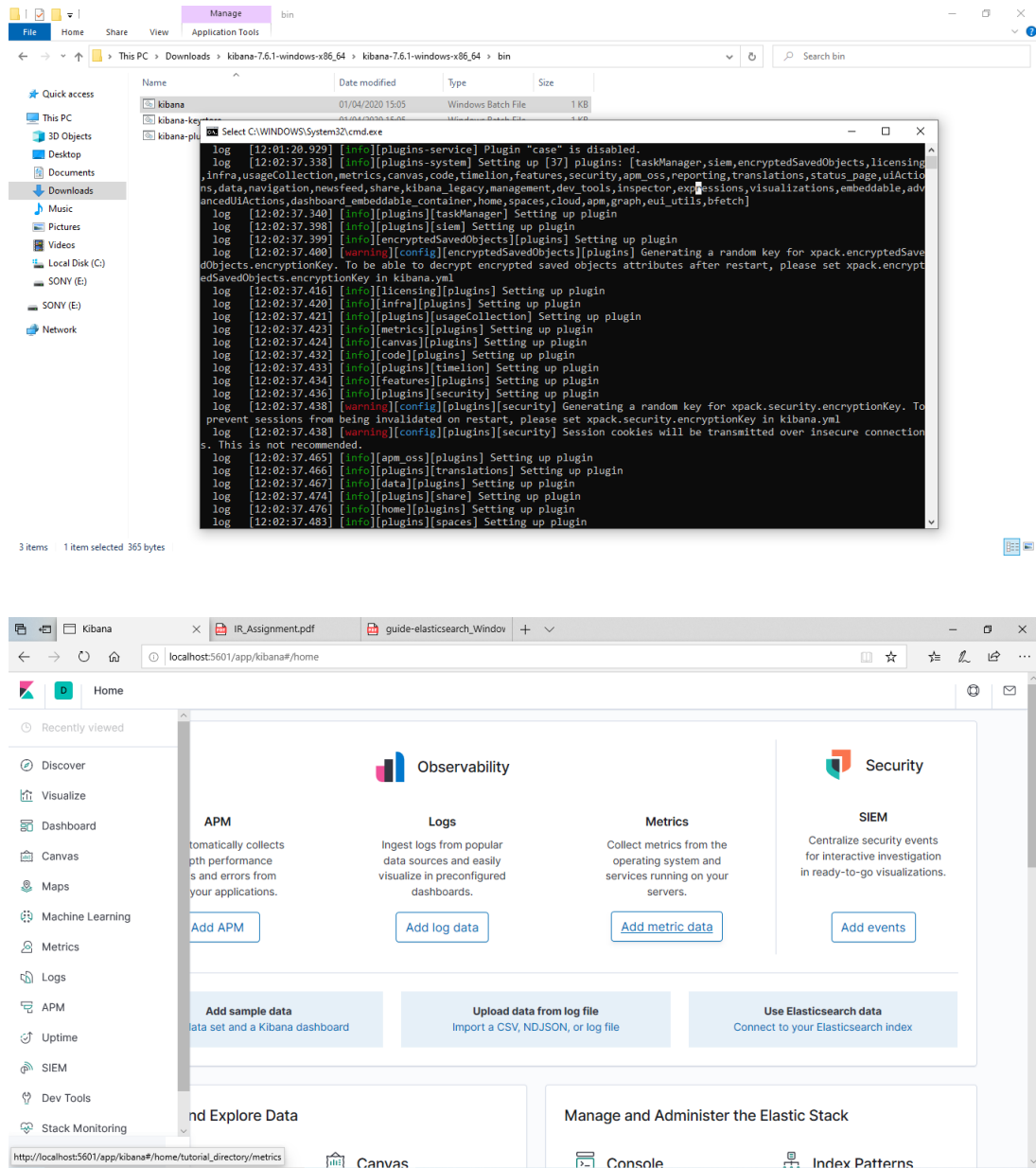


Figure 2: Running Kibana

- Python Python 3.8.0 to access the data in Kibana for indexing and searching.

Three packages of python have been used i.e elasticsearch to connect with Elasticsearch software running in the system, json_lines to read lines where each line is json encoded, json for reading and writing from json file and to convert string into json file. Along with these three important packages we used mathematical library numpy and nltk library.

1.2 Indexing

The group of documents having indistinguishable attributes referred to as index. It is recognised by its distinctive name while carrying out task such as searching. Cluster are nodes that collectively defines the data and are capable enough to provide indexing and searching of documents. There is no limit of indexes with in an cluster it can be as many as we want. Similarly, there is no limitation of storing number of documents in an index [2].

From the dataset, we are using 8000 articles for which we are going to create index using Elasticsearch and python in *indexing()* function of “searchengine.py” python file by following certain steps.

- Download the file used for indexing [signalmedia-1m.jsonl](#).
- Read this json file by using json_lines package.
- Output obtained in the form of string.
- Convert string into json form using json.dumps() function of json package.
- Load the file in Elasticsearch.

To see the index created in Kibana after proper execution of file.

- Open Kibana site from the browser.
- Click on Management option on left side of the software.
- Select Index Management.
- Check whether the index is created or not.
- Click on the respective index and access mapping tab which are automatically generated while feeding the data in Elasticsearch.
- Fields of data: id, content, title, media-type, source, and published

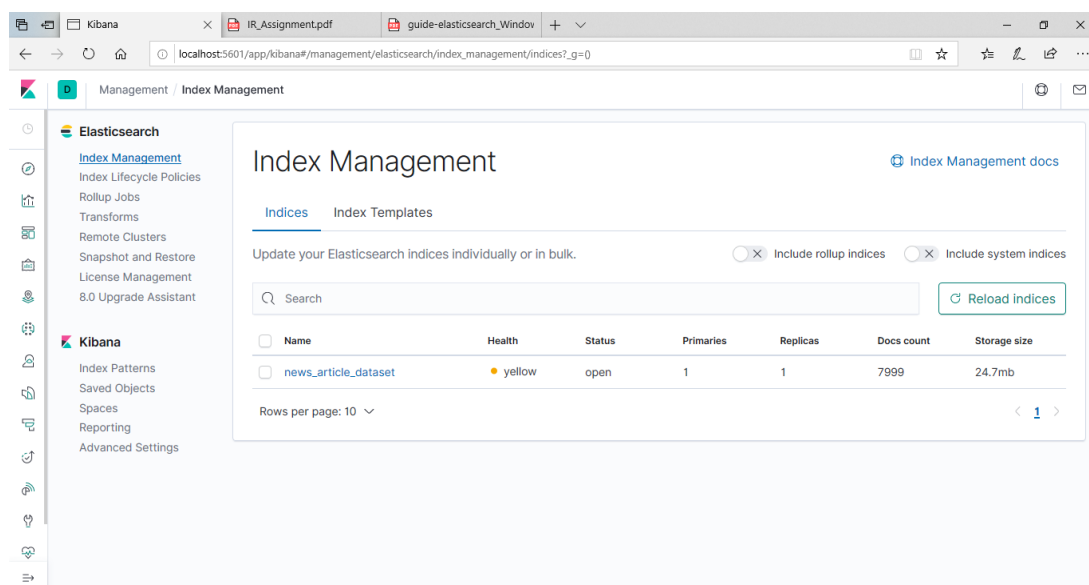
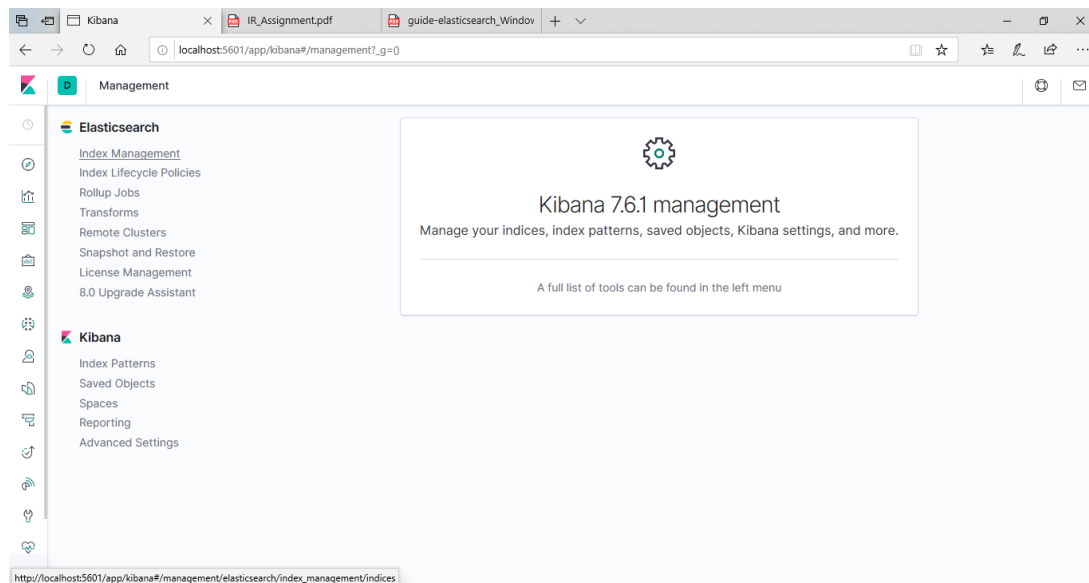


Figure 3: Indexing Json Dataset

In order to see the mapping of the dataset, we need to click on the index first then click to mapping. The mapping are created automatically after we load the dataset onto elastic search using the *indexing()* function.

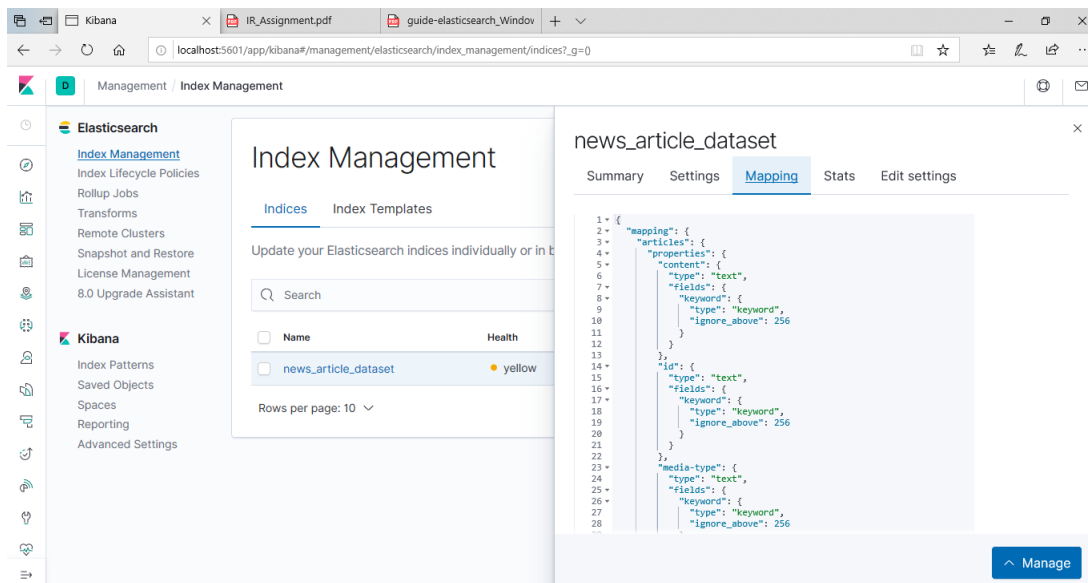


Figure 3: Mapping Json Dataset

1.3 Searching

In order to use Elasticsearch there is a need of interface. So we are using Kibana which is a web interface required to do the manipulation as well as visualization of data to execute more complex queries [3]. Along with this, we are searching uploaded files with the help of created index in Kibana site.

- Open Kibana site from the browser.
- Click on Management option on bottom left side of the software.
- Select Index Patterns.
- Click on Create Index Patterns.

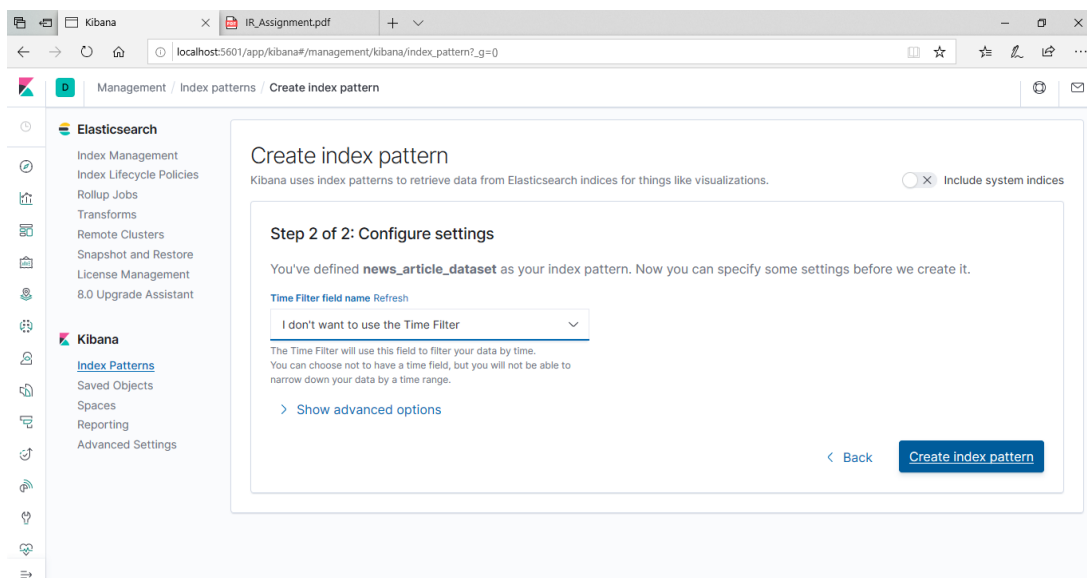
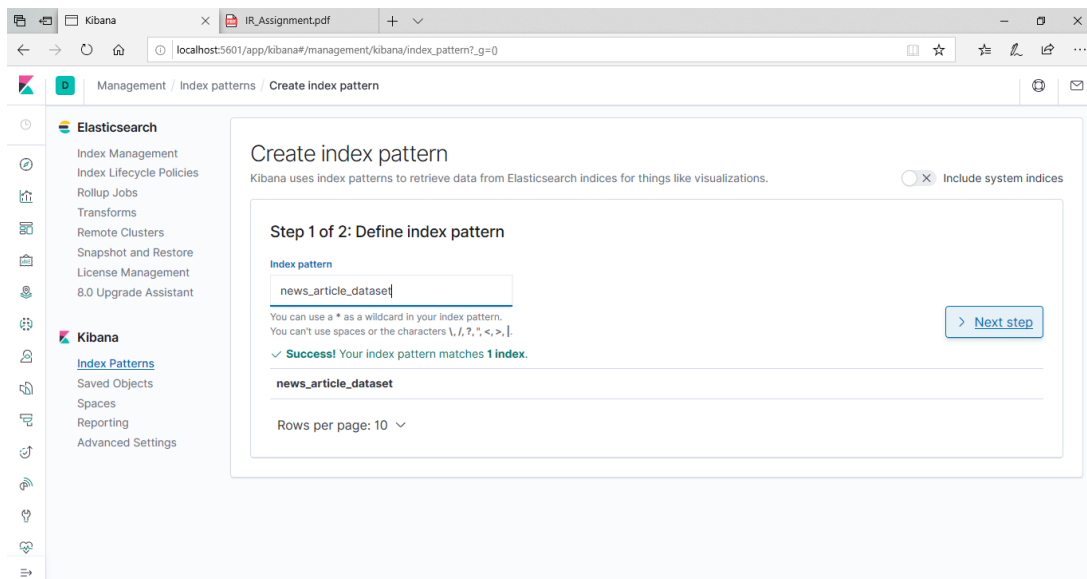


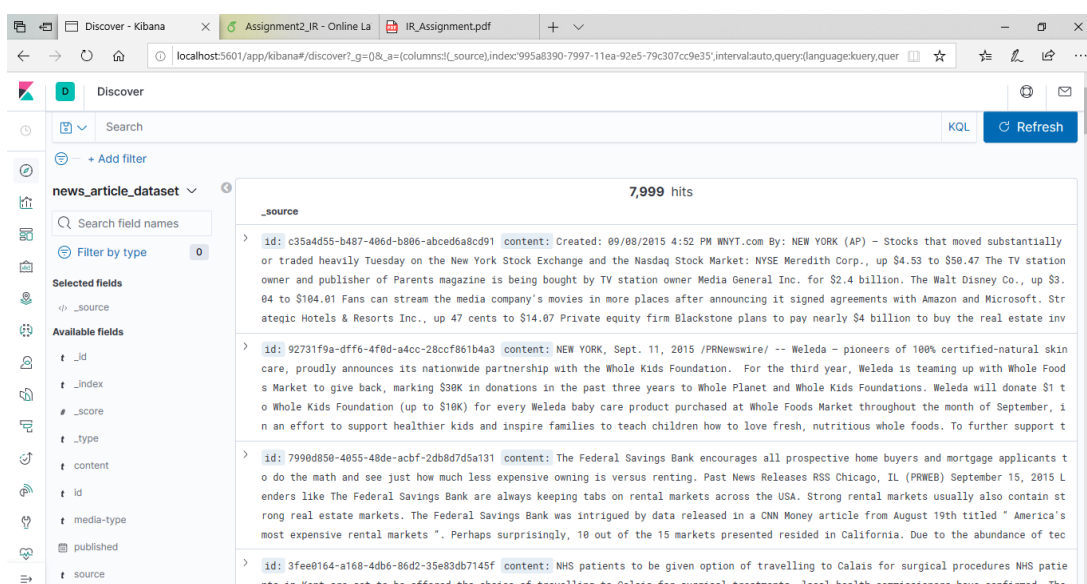
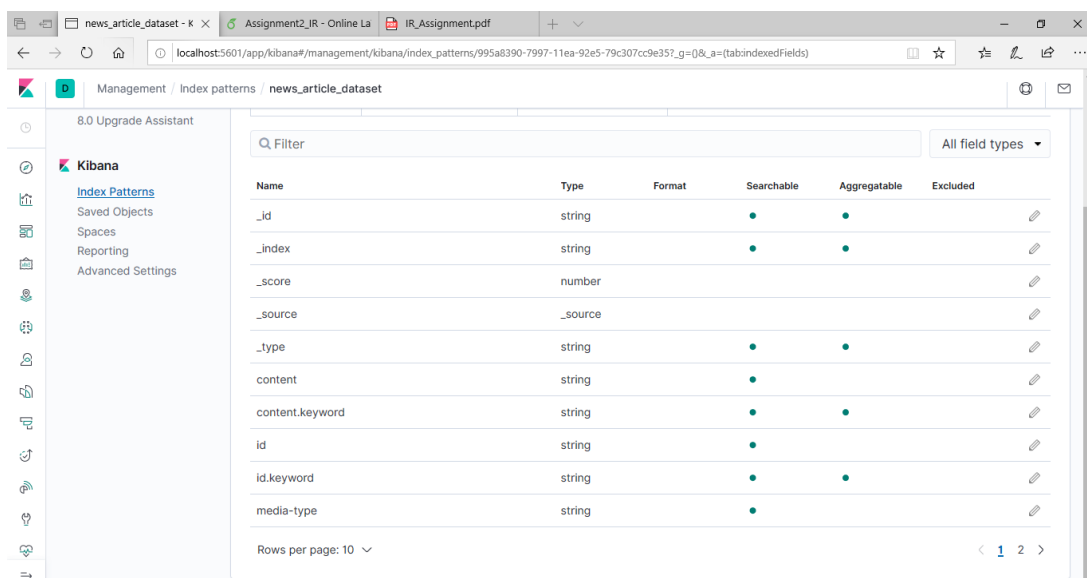
Figure 4: Creating Index Pattern

After creating indexes, we need to search indexes which can be done by using discover tab which is visible on upper left corner of Kibana site.

Searching operation is carried out by following below steps:

- Select the keyword to search just like media-type.
- Go to edit filter in Kibana and add the filter with the selected keywords.
- Save the search.

- Results of search are visible.
- We can also add multiple filters just like accidents.
- if we search for accident's then the action related to accidents are visible. It indicates that, there are articles about accident or occurrence of unfortunate incidents which contain in these keywords
- In case we are specifying search and title then consequently it will give the output in a precise way.



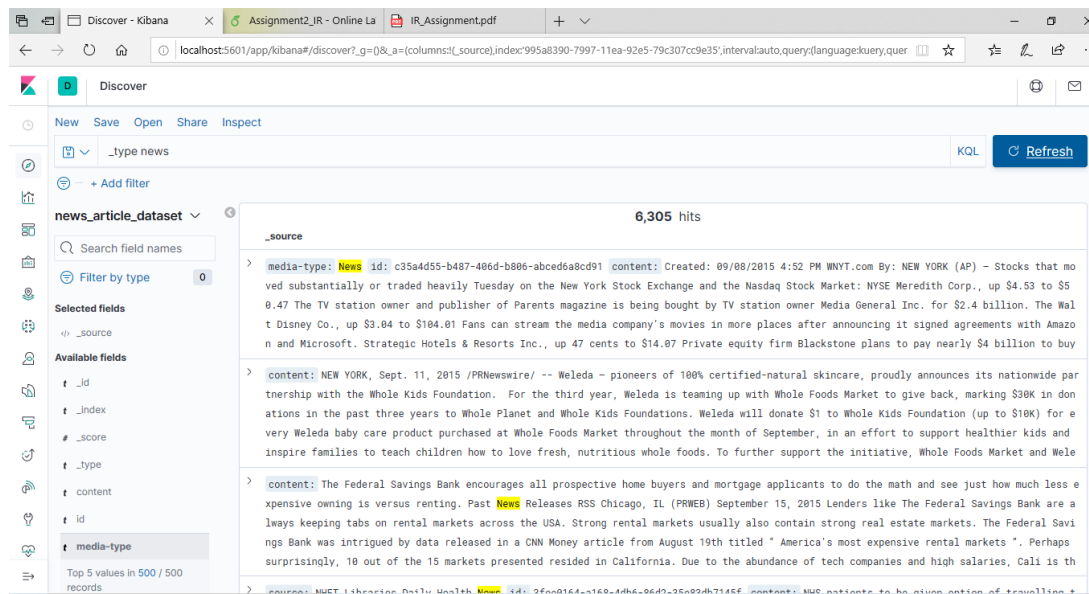


Figure 5: Searching Information

1.4 Building a Test Collection

Test collection are developed in order to test the efficiency of information retrieval systems. It gives us the opportunity to explore how can we optimize our search engine for the particular collection we are indexing and make it more efficient. Following test collection has been devised.

Outputs of 6 different queries:

1. Query 1: Query will search title and content of the article.

```
{
  "query": {
    "bool": {
      "should": [
        {
          "match": {
            "title": "title"
          }
        },
        {
          "match": {
            "content": "content"
          }
        }
      ]
    }
  },
  "size": 8000
}
```

2. Query matches with the exact title content type

```

    {"query":
      {"match_phrase":
        {"title":title
        }
      }
    },size=8000

```

3. Query for content which have higher importance than title

```

{"query":
  {"multi_match":
    {"query":keyword,"fields":
      ["content^2", "title"
      ]
    }
  }
},size=8000

```

4. Query for content to begin with letter

```

{"query":
  {"wildcard":
    {"content": content
    }
  }
},size=8000

```

5. Query for Source and Content

```

{"query":
  {"bool":
    {"should":
      [
        {"match":
          {"source":source
          }
        },
        {"match":
          {"content":content
          }
        }
      ]
    }
  }
},size=8000

```

6. Query to search Media type and content range of published date

```

{"query":
  {"bool":
    {"should":
      [
        {"match":
          {"media-type":media
        }
      },
      {"match":
        {"content":content
      }
    },
    {"range":
      {"published":
        {"gte": start,"lte": end,"format":
          "yyyy/MM/dd||yyyy"
        }
      }
    ]
  }
},size=8000

```

1.5 Evaluation

Evaluation are also used primarily for measuring the effectiveness of system. It is carried out basically to find out goodness of the system and its performance.

For the evaluation, we are using mean average precision popular metrics as the measure for accuracy. Mean average is the combination of recall and precision for the retrieved outcomes. It is calculated as ratio of number of queries in a set divided by average precision for a given query [4].

The *mapr()* function has been used to carry out the evaluation.

Follow the below steps to evaluate:

- Initialization of precision and recall values equal to zero.
- Evaluate ranked recall and precision for the search engine.
- Array data structure is used to store indexes and sort them for every document retrieved which helps us to compute precision and recall values.
- It is a cycle where precision and recall are calculated at each stage continuously.
- After calculating precision and recall for each document,compute the average of all values of the precision and recall.
- Output is shown with the respective document id,precision,recall,mean average preci-

sion,mean average recall.

```
Python 3.8.0 Shell
File Edit Shell Debug Options Window Help
7964 0.00012556504269211453 0.000125 500
7965 0.00012554927899168997 0.000125 500
7966 0.00012553351746915893 0.000125 500
7967 0.000125517760763148 0.000125 500
7968 0.0001255020080321285 0.000125 500
7969 0.00012549625925461163 0.000125 500
7970 0.00012547051442910815 0.000125 500
7971 0.00012545477355413375 0.000125 500
7972 0.0001254390366281987 0.000125 500
7973 0.00012542330364981815 0.000125 500
7974 0.00012540757461750688 0.000125 500
7975 0.00012539184952978957 0.000125 500
7976 0.00012537612838515847 0.000125 500
7977 0.00012536041118214867 0.000125 500
7978 0.000125344697919278 0.000125 500
7979 0.00012532898859506203 0.000125 500
7980 0.00012531328320602005 0.000125 500
7981 0.0001252975817564721 0.000125 500
7982 0.00012528188423953897 0.000125 500
7983 0.0001252661906551422 0.000125 500
7984 0.000125250501002004 0.000125 500
7985 0.000125234801527864746 0.000125 500
7986 0.00012521913348359628 0.000125 500
7987 0.000125203455615375 0.000125 500
7988 0.00012518778167250875 0.000125 500
7989 0.00012517211163352359 0.000125 500
7990 0.0001251564555694618 0.000125 500
7991 0.00012514078338130396 0.000125 500
7992 0.00012512512512512512 0.000125 500
7993 0.00012510947078693858 0.000125 500
7994 0.00012509382036527395 0.000125 500
7995 0.00012507817385966166 0.000125 500
7996 0.0001250625312656328 0.000125 500
7997 0.00012504689258471927 0.000125 500
7998 0.00012503125781445363 0.000125 500
7999 0.0001250156269336918 0.000125 500
8000 0.000125 0.000125 500
8001 0.00012498437685288088 0.000125 500
Mean Average Precision: 0.00034647206733098287
Mean Average Recall: 0.00011720375000001313
Would you like to query the search engine? press y/n|
```

Figure 7: Evaluation

1.6 Engineering a Complete System

The system has been developed using Elasticsearch, Kibana and python. The system initiate by taking the confirmation that the file is indexed. Then further proceed by firing the query and results will be displayed, followed by displaying precision and recall at every stage of the system and at last showing mean average precision.

```
Python 3.8.0 Shell
File Edit Shell Debug Options Window Help
Python 3.8.0 (tags/v3.8.0:fa519fd, Oct 14 2019, 19:37:50) [MSC v.1916 64 bit (AMD64)] on win32
Type "help", "copyright", "credits" or "license()" for more information.
>>>
= RESTART: C:\Users\USER\Downloads\Information-Retrieval-master\searchengine.py
Is the file Indexed in Elastic Search? yes
Would you like to query the search engine? press y/n:y

1.Searching for word through title ranged using published date:
2.Title and content is asked:
3.Media type and content range of published date:
4.Any title but content should contain start phrase:
5.Title should be exact:
6.Content have higher importance than title:
7.Content should begin with a letter:
8.Source and content is asked:
9.Enter 2 keywords, will return if both words appear together in content:
10.Enter 2 keywords, will return if both words appear together or seperately in content:
Enter a choice of query:5
Enter the Exact Title:10 Great Game Soundtracks for Studying
Document Index is: 500
Document Search Score is: 39.642914
Document Media Type: Blog
Document Title: 10 Great Game Soundtracks for Studying
Source: IGN All
Published on: 2015-09-03T18:55:39Z
Document Content: School is back in session, which means inevitable late nights putting the finishing touches on an essay, or cramming as many formulas into your head as you can before the next exam. Since most game music is designed to create a backdrop for gameplay and help focus the player on solving whatever problems lie before them, it's no coincidence that many of the best game soundtracks are also great resources for studying. Here are 10 of the best you can download right now:
```

Figure 7: Engineering a Complete System

1.7 *Crowdsourcing*

Crowd sourcing is the process where large group of people contribute by investing their time for a common objective may be any research or any other goal such as innovation, problem-solving. The exercise we have completed majorly test our memory power. It majorly tests our capability to remember the videos and analysing the patterns of different brains. It is a simple test to learn the efficiency of our mind. It is based on reflexive actions taken by the user by coordinating their brain with their action by clicking space bar whenever we found repeated videos. It increases our concentration by doing some task with full dedication for 20 minutes time span. The user experience can be enhanced by having proper instructions at what time we should press space bar at the end of the video or in the middle. The interface can be more lively for the exercise. It can be made more interesting by adding the idea of ramification to it because in some sense it resembles towards playing any mindcentric game. After some videos of seven second there is white blank screen appears which makes the experience somewhat less interesting. We have learned to be more focused and attentive towards our task and user experience can be improved by making the platform more interactive and attractive.

2 *References*

- [1] <https://towardsdatascience.com/an-overview-on-elasticsearch-and-its-usage-e26df1d1d24a>
- [2] https://moodle.essex.ac.uk/pluginfile.php/1007586/course/section/139844/assignment_2_306_706.pdf?time=1588800000
- [3] <https://medium.com/@victorsmelopoa/an-introduction-to-elasticsearch-with-kibana-78071db3704>
- [4] <https://towardsdatascience.com/breaking-down-mean-average-precision-map-ae462f623a52>