

# CE706 - Information Retrieval 2021

## Assignment 1

2004532

### Instructions for running your system (Engineering a Complete System)

#### Elastic Search:

- Download Elastic Search for Windows

[https://artifacts.elastic.co/downloads/elasticsearch/elasticsearch-7.11.1-windows-x86\\_64.zip](https://artifacts.elastic.co/downloads/elasticsearch/elasticsearch-7.11.1-windows-x86_64.zip)

- Unzip the downloaded file
- Navigate to bin folder and run elasticsearch.bat
- Check <http://localhost:9200/> to see elastic search running.

The screenshot shows a REST client interface with a GET request to `http://localhost:9200` sent successfully. The response is a JSON object representing the Elasticsearch status.

KEY	VALUE	DESCRIPTION
Key	Value	Description

Body: Cookies Headers (3) Test Results Status: 200 OK Time: 16 ms Size: 431 B Save Response

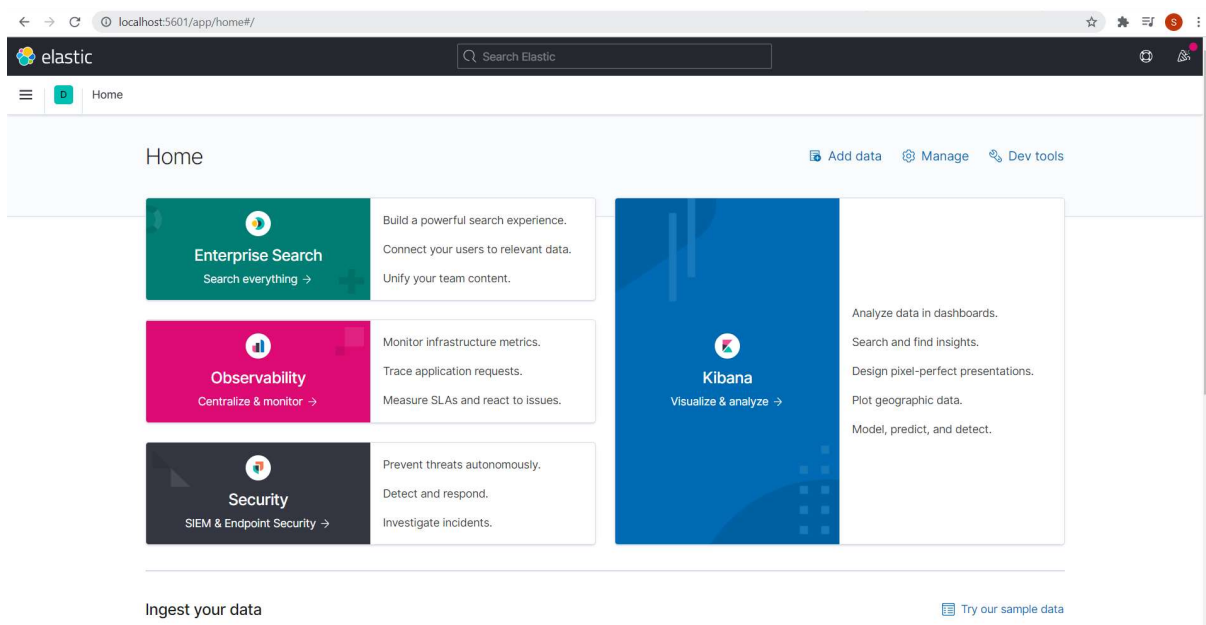
```
1 {
2   "name": "MSI",
3   "cluster_name": "elasticsearch",
4   "cluster_uuid": "D7K5ddpMRgCdLoxXIrpfg",
5   "version": {
6     "number": "7.11.0",
7     "build_flavor": "default",
8     "build_type": "zip",
9     "build_hash": "8ced7813d6f16d2ef30792e2fcde3e755795ee04",
10    "build_date": "2021-02-08T22:44:01.320463Z",
11    "build_snapshot": false,
12    "lucene_version": "8.7.0",
13    "minimum_wire_compatibility_version": "6.8.0",
14    "minimum_index_compatibility_version": "6.0.0-beta1"
15  },
16   "tagline": "You Know, for Search"
17 }
```

#### Kibana:

- Download Kibana for Windows3

[https://artifacts.elastic.co/downloads/kibana/kibana-7.11.1-windows-x86\\_64.zip](https://artifacts.elastic.co/downloads/kibana/kibana-7.11.1-windows-x86_64.zip)

- Unzip the downloaded file
- Navigate to bin folder and run kibana.bat
- Check <http://localhost:5601/> to see kibana running.



### Python Jupyter and Elastic Search:

- Make sure you have python and jupyter correctly set up.
- Install the following python libraries by running these commands:

```
pip install -U scikit-learn
```

```
pip install -U pip setuptools wheel
```

```
pip install -U spacy
```

```
python -m spacy download en_core_web_sm
```

```
pip install elasticsearch
```

## Indexing

An index can be thought of as an optimized collection of documents and each document is a collection of fields, which are the key-value pairs that contain your data. In this assignment I have set the index to value “covid” because we are dealing with covid articles. Indexing allows faster search of documents. An index consists of documents, and a document consists of fields.

I downloaded the dataset for covid articles i.e metadata.csv from Kaggle site. The dataset is in csv form. So before uploading the data to elastic search, I added two new fields to the dataset i.e. index and id. Specifying index allows us to dump data to elasticsearch to specific index. It is not necessary to add id field to the data because otherwise elasticsearch automatically generates an id for each document. However in the dataset there is a field called cord\_uid so I set this as an id for each document. Moreover I found out that there are columns (mag\_id, arxiv\_id) that contain almost no data so I dropped these two columns from the dataset.

```
df.isnull().sum()
```

```
cord_uid      0
sha           299973
source_x      0
title         223
doi           199517
pmcid         290607
pubmed_id     231259
license       0
abstract      124997
publish_time   218
authors       13110
journal       30185
mag_id        450385
who_covidence_id 268359
arxiv_id      444490
pdf_json_files 299973
pmc_json_files 329971
url           180409
s2_id         41014
dtype: int64
```

Some of the values in title, abstract and authors column are missing so I excluded those documents.

```
df[pd.isna(df.title) | pd.isna(df.abstract) | pd.isna(df.authors)].shape
```

```
(128062, 17)
```

```
df=df[pd.notna(df.title) & pd.notna(df.abstract) & pd.notna(df.authors)]
```

Some of the values in publish\_time field were also missing so we removed those records from the dataset as well.

```
df.publish_time=pd.to_datetime(df.publish_time)
```

```
df[pd.isna(df.publish_time)].shape
```

```
(188, 17)
```

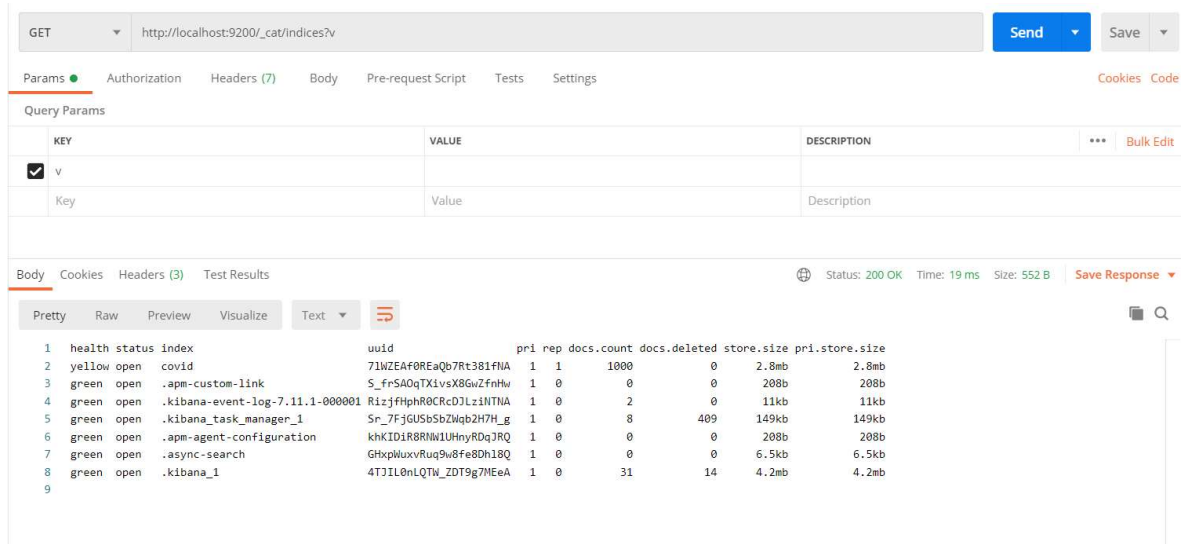
```
df=df[pd.notna(df.publish_time)]
```

After that I selected 1000 documents from the dataset and uploaded it to elasticsearch. Moreover I specified mapping for two fields i.e. authors and keywords in data as keyword type in elastic search. I created keywords field in data that contains keywords from textual data in the dataset. Specifying these fields as keyword type allows faster search when we query for these respective fields. I also specified publish\_time field as date type in elastic search.

```
from elasticsearch import helpers
from elasticsearch import Elasticsearch

es=Elasticsearch()
es.indices.create(index="covid",ignore=400, body={
    "mappings": {
        "properties": {
            "authors":{
                "type": "keyword",
            },
            "keywords":{
                "type": "keyword",
            },
            "publish_time":{
                "type": "date",
                "format": "yyyy-MM-dd"
            }
        }
    }
})
```

After uploading and indexing the data we could verify it on postman:



The screenshot shows a Postman interface with a GET request to `http://localhost:9200/_cat/indices?v`. The response is a table of Elasticsearch indices. The table has columns: `health`, `status`, `index`, `uuid`, `pri`, `rep`, `docs.count`, `docs.deleted`, `store.size`, `pri.store.size`.

health	status	index	uuid	pri	rep	docs.count	docs.deleted	store.size	pri.store.size
yellow	open	covid	71WZEAf0REaQb7Rt381fNA	1	1	1000	0	2.8mb	2.8mb
green	open	.apm-custom-link	S_frSAQqTXivsX8GwZfnHw	1	0	0	0	208b	208b
green	open	.kibana-event-log-7.11.1-000001	RizjfhPhR0CRcDjLzINTNA	1	0	2	0	11kb	11kb
green	open	.kibana_task_manager_1	Sr_7FjGUSbSbZwqb2H7H_g	1	0	8	409	149kb	149kb
green	open	.apm-agent-configuration	khKIDiR8RNNiUHnyRDqJRQ	1	0	0	0	208b	208b
green	open	.async-search	GhxpWuxvRuq9w8fe8Dh18Q	1	0	0	0	6.5kb	6.5kb
green	open	.kibana_1	4T3JL8nLQTW_ZDT9g7MEaA	1	0	31	14	4.2mb	4.2mb

## Sentence Splitting, Tokenization and Normalization

In the dataset, there are three columns that contain textual data i.e. title, abstract and authors. For authors field, the name of authors are separated by “;”, so I split the text in authors field with “;” delimiter and stored it as an array. In mapping I specified the author field as keyword type so that it would be easier to search for specific author publications. For title and abstract, I merged the content from these two fields to a single field. In the initial preprocessing I removed extra spaces, new line/tab characters, punctuations and numbers/digits from the text using python and regex expressions. After that I used spacy library in python to tokenize the text and for removing stop words.

Before:

Debate: Transfusing to normal haemoglobin levels will not improve outcome Recent evidence suggests that critically ill patients are able to tolerate lower levels of haemoglobin than was previously believed. It is our goal to show that transfusing to a level of 100 g/l does not improve mortality and other clinically important outcomes in a critical care setting. Although many questions remain, many laboratory and clinical studies, including a recent randomized controlled trial (RCT), have established that transfusing to normal haemoglobin concentrations does not improve organ failure and mortality in the critically ill patient. In addition, a restrictive transfusion strategy will reduce exposure to allogeneic transfusions, result in more efficient use of red blood cells (RBCs), save blood overall, and decrease health care costs.

After initial preprocessing:

Debate Transfusing to normal haemoglobin levels will not improve outcome Recent evidence suggests that critically ill patients are able to tolerate lower levels of haemoglobin than was previously believed It is our goal to show that transfusing to a level of g l does not improve mortality and other clinically important outcomes in a critical care setting Although many questions remain many laboratory and clinical studies including a recent randomized controlled trial RCT have established that transfusing to normal haemoglobin concentrations does not improve organ failure and mortality in the critically ill patient In addition a restrictive transfusion strategy will reduce exposure to allogeneic transfusions result in more efficient use of red blood cells RBCs save blood overall and decrease health care costs

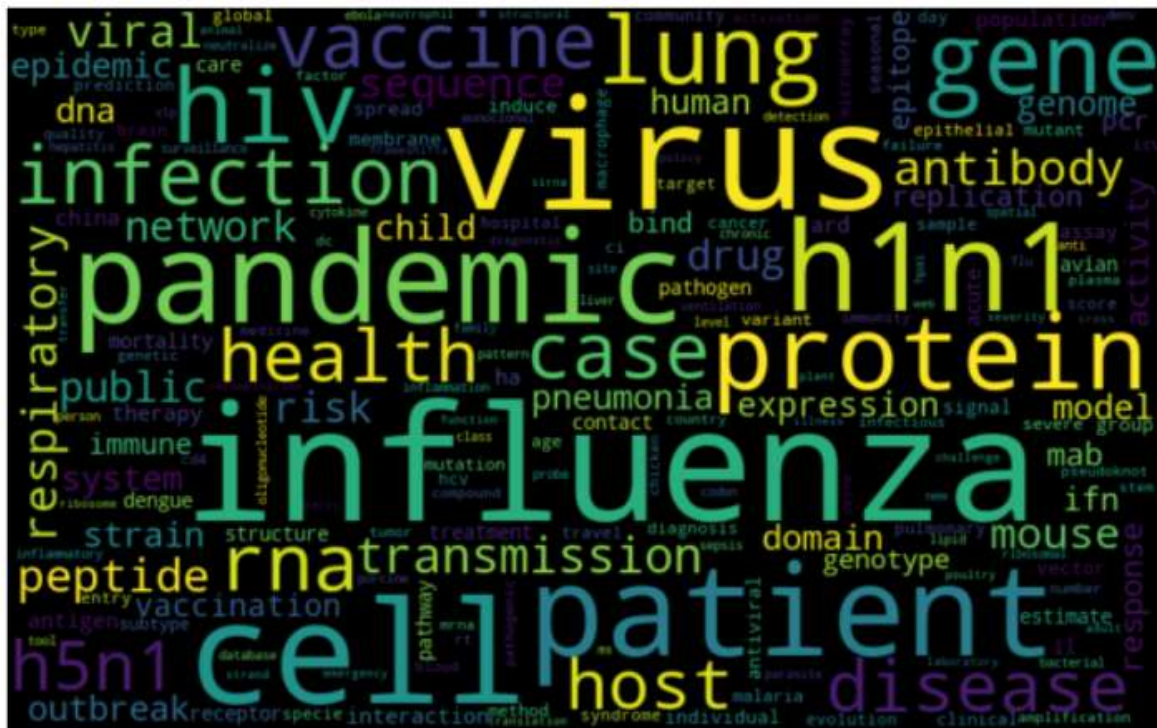
After spaCy preprocessing:

```
['debate', 'transfuse', 'normal', 'haemoglobin', 'level', 'improve', 'outcome', 'recent', 'evidence', 'suggest', 'critically', 'ill', 'patient', 'able', 'tolerate', 'low', 'level', 'haemoglobin', 'previously', 'believe', 'goal', 'transfusing', 'level', 'g', 'l', 'improve', 'mortality', 'clinically', 'important', 'outcome', 'critical', 'care', 'setting', 'question', 'remain', 'laboratory', 'clinical', 'study', 'include', 'recent', 'randomize', 'controlled', 'trial', 'rct', 'establish', 'transfuse', 'normal', 'haemoglobin', 'concentration', 'improve', 'organ', 'failure', 'mortality', 'critically', 'ill', 'patient', 'addition', 'restrictive', 'transfusion', 'strategy', 'reduce', 'exposure', 'allogeneic', 'transfusion', 'result', 'efficient', 'use', 'red', 'blood', 'cell', 'rbc', 'save', 'blood', 'overall', 'decrease', 'health', 'care', 'cost']
```



## Selecting Keywords

After preprocessing the text data, I removed all the unnecessary text from data and only left with words that can be keyword. For selecting keyword I used tf-idf (term frequency-inverse document frequency). TF-IDF allows us to compute a score against each word in a document that shows how important a word is in a document in a collection of documents. Lower tf-idf score shows that the word is more common in different documents and high tf-idf score shows that a word can be a keyword of document. I only selected top 10 keywords for each document and stored it with respective document in the dataset. The wordcloud for our keywords in all documents is shown below:



## Stemming or Morphological Analysis

Stemming and Lemmatization both generate the root form of the words. Moreover stemming might produce words that are not in English vocabulary but with lemmatization we don't need to worry about that. spaCy doesn't contain any function for stemming so for morphological analysis I used lemmatization that is supported by spaCy and works very well. Here are few examples of lemmatization in our dataset:

Before Lemmatization

[ 'Managing', 'emerging', 'infectious', 'diseases', 'federal', 'system', 'impediment', 'effective', 'laws', 's', 's', 'HIV', 'AI DS', 'emerging', 'infectious', 'disease', '2004', 'saw', 'emergence', 'SARS', 'Avian', 'influenza', 'Anthrax', 'man', 'form', 'bioterrorism', 'Emergency', 'powers', 'legislation', 'Australia', 'patchwork', 'Commonwealth', 'quarantine', 'laws', 'State', 'Territory', 'based', 'emergency', 'powers', 'public', 'health', 'legislation', 'time', 'review', 'legislation', 'time', 'consi deration', 'efficacy', 'legislation', 'country', 'wide', 'perspective', 'age', 'consider', 'possibility', 'mass', 'outbreaks', 'communicable', 'diseases', 'ignore', 'jurisdictional', 'boundaries' ]

After Lemmatization

```
[ 'management', 'emerge', 'infectious', 'disease', 'federal', 'system', 'impediment', 'effective', 'law', 's', 's', 'hiv', 'aids', 'emerge', 'infectious', 'disease', '2004', 'see', 'emergence', 'sars', 'avian', 'influenza', 'anthrax', 'man', 'form', 'bioterrorism', 'emergency', 'power', 'legislation', 'australia', 'patchwork', 'commonwealth', 'quarantine', 'law', 'state', 'territory', 'base', 'emergency', 'power', 'public', 'health', 'legislation', 'time', 'review', 'legislation', 'time', 'consideration', 'efficacy', 'legislation', 'country', 'wide', 'perspective', 'age', 'consider', 'possibility', 'mass', 'outbreak', 'communicable', 'disease', 'ignore', 'jurisdictional', 'boundary']
```

# Searching

I used Kibana for searching the indexed documents that we uploaded to elastic search. Kibana gives a gui for better visualization and selection of fields in our data.

elastic

Search Elastic

Stack Management / Index patterns

Ingest

Ingest Node Pipelines

Data

Index Management

Index Lifecycle Policies

Snapshot and Restore

Rollup Jobs

Transforms

Remote Clusters

Alerts and Insights

Alerts and Actions

Reporting

Kibana

Index Patterns

Saved Objects

Tags

Spaces

Advanced Settings

Stack

You have data in Elasticsearch.

Now, create an index pattern.

Kibana requires an index pattern to identify which indices you want to explore. An index pattern can point to a specific index, for example, your log data from yesterday, or all indices that contain your log data.

Create index pattern

Want to learn more? [Read documentation](#)

elastic

Search Elastic

Stack Management / Index patterns / covid\*

Ingest

Ingest Node Pipelines

Data

Index Management

Index Lifecycle Policies

Snapshot and Restore

Rollup Jobs

Transforms

Remote Clusters

Alerts and Insights

Alerts and Actions

Reporting

Kibana

Index Patterns

Saved Objects

Tags

Spaces

Advanced Settings

Stack

covid\*

This page lists every field in the **covid\*** index and the field's associated core type as recorded by Elasticsearch. To change a field type, use the Elasticsearch [Mapping API](#)

Fields (34) | Scripted fields (0) | Field filters (0)

Search

All field types

Name	Type	Format	Searchable	Aggregatable	Excluded
_id	string		•	•	
_index	string		•	•	
_score	number				
_source	_source				
_type	string		•	•	
abstract	string		•		
abstract.keyword	string		•	•	
authors	string		•	•	

The search results for different fields are given below:

elastic Search Elastic

Discover

Search

1,000 hits

Filter by type 0

Available fields

- \_id
- \_index
- \_score
- \_type
- abstract
- authors
- cord\_uid
- doi
- journal
- keywords
- license
- pdf\_json\_files
- pmc\_json\_files
- pmcid

cord\_uid: 13227806 sha: 7d624acf891472c3452b1deb89e8dc44a48d8c94 source\_x: PMC title: GIDEON: a comprehensive Web-based resource for geographic medicine doi: 10.1186/1476-072x-4-10 pmcid: PMC1090610 pubmed\_id: 15847698 license: cc-by abstract: GIDEON (Global Infectious Diseases and Epidemiology Network) is a web-based computer program designed for decision support and informatics in the field of Geographic Medicine. The first of four interactive modules generates a ranked differential diagnosis based on patient signs, symptoms, exposure history and country of disease acquisition. Additional options include syndromic disease surveillance capability and simulation of bioterrorism scenarios. The second module accesses detailed and current information regarding the status of

cord\_uid: yz2wbpuu sha: d786e82e2ebc8571e70654d1acf0ffdd3969ab2d3 source\_x: PMC title: Globalization and Health doi: 10.1186/1744-8603-1-1 pmcid: PMC1143779 pubmed\_id: 15847699 license: cc-by abstract: This debut editorial of Globalization and Health introduces the journal, briefly delineating its goals and objectives and outlines its scope of subject matter. 'Open Access' publishing is expected to become an increasingly important format for peer reviewed academic journals and that Globalization and Health is 'Open Access' is appropriate. The rationale behind starting a journal dedicated to globalization and health is three fold: Firstly: Globalization is reshaping the social geography within which we might strive to create health or prevent

cord\_uid: kvhoa2se sha: ada57a9f084a016deb20b6f83c442072e637955 source\_x: PMC title: The 'polysemous' codon--a codon with multiple amino acid assignment caused by dual specificity of tRNA identity. doi: 10.1093/emboj/16.5.1122 pmcid: PMC1169711 pubmed\_id: 9118950 license: no-cc abstract: In some Candida species, the universal CUG leucine codon is translated as serine. However, in most cases, the serine tRNAs responsible for this non-universal decoding (tRNA(Ser)CAG) accept in vitro not only serine, but also, to some extent, leucine. Nucleotide replacement experiments indicated that m1G37 is critical for leucylation activity. This finding was supported by the fact that the tRNA(Ser)CAGs possessing the leucylation activity always have m1G37, whereas that of

cord\_uid: cg134ykt sha: 6eeaffdb22b4e310867262d1891a1442769c7803 source\_x: PMC title: A universal BMV-based RNA recombination system-how to search for general rules in RNA recombination doi: 10.1093/nar/gni106 pmcid: PMC1174899 pubmed\_id: 16802784 license: no-cc abstract: At present, there is no doubt that RNA recombination is one of the major factors responsible for the generation of new RNA viruses and retroviruses. Numerous experimental systems have been created to investigate this complex phenomenon. Consequently, specific RNA structural motifs mediating recombination have been identified in several viruses.

elastic Search Elastic

Discover

Search

1 hit

Filter by type 0

Available fields

- \_id
- \_index
- \_score
- \_type
- abstract
- authors
- cord\_uid
- doi
- journal
- keywords
- license
- pdf\_json\_files
- pmc\_json\_files
- pmcid

cord\_uid: kvhoa2se sha: ada57a9f084a016deb20b6f83c442072e637955 source\_x: PMC title: The 'polysemous' codon--a codon with multiple amino acid assignment caused by dual specificity of tRNA identity. doi: 10.1093/emboj/16.5.1122 pmcid: PMC1169711 pubmed\_id: 9118950 license: no-cc abstract: In some Candida species, the universal CUG leucine codon is translated as serine. However, in most cases, the serine tRNAs responsible for this non-universal decoding (tRNA(Ser)CAG) accept in vitro not only serine, but also, to some extent, leucine. Nucleotide replacement experiments indicated that m1G37 is critical for leucylation activity. This finding was supported by the fact that the tRNA(Ser)CAGs possessing the leucylation activity always have m1G37, whereas that of





