

Text mining for exploration of COVID-19 severity factors

Khang Duy LAI - Mariia KLIMINA

University Paris Cité

UFR des Sciences Fondamentales et Biomédicales

May 12, 2022

2022-05-12

Text mining for exploration of COVID-19 severity factors

1 Introduction

2 State of the art

3 Data exploration

4 Data preprocessing

5 Data processing

6 Result

2022-05-12

Text mining for exploration of COVID-19 severity factors

- 1 Introduction
- 2 State of the art
- 3 Data exploration
- 4 Data preprocessing
- 5 Data processing
- 6 Result

Introduction●	State of the art○	Data exploration○○○○	Data preprocessing○○	Data processing○○○○○○○	Result○○○○
Introduction					
COVID-19 is the disease caused by the Sar-COV-2 virus that originated in China at the end of the year 2019. Over the time, studies have shown that there is some form of background diseases and risk factors that can hugely affect the severity cases rate of COVID-19. This project will apply NLP and text mining methods in order to explore the CORD-19 dataset and extract background diseases and risk factors.					
Khang Duy LAI - Mariia KLIMINA					
Text mining for exploration of COVID-19 severity factors					

2022-05-12	Text mining for exploration of COVID-19 severity factors	
	└ Introduction	
	└ Introduction	
Introduction		
COVID-19 is the disease caused by the Sar-COV-2 virus that originated in China at the end of the year 2019. Over the time, studies have shown that there is some form of background diseases and risk factors that can hugely affect the severity cases rate of COVID-19. This project will apply NLP and text mining methods in order to explore the CORD-19 dataset and extract background diseases and risk factors.		

State of the art

In this project we used multiple state of the art NLP and Data Science libraries.

- Numpy,Pandas: Formatting the data and the calculations.
- Matplotlib: Library for drawing the charts and figures.
- Scikit-learn: LDA and T-SNE models.
- Spacy,Gensim, and NLTK: Important NLP libraries.
- Scispacy: NER,Spacy models for science papers.
- Bokeh: A library for visualising interacted charts.

2022-05-12

Text mining for exploration of COVID-19 severity factors

- └ State of the art
- └ State of the art

State of the art

In this project we used multiple state of the art NLP and Data Science libraries.

- Numpy,Pandas: Formatting the data and the calculations.
- Matplotlib: Library for drawing the charts and figures.
- Scikit-learn: LDA and T-SNE models.
- Spacy,Gensim, and NLTK: Important NLP libraries.
- Scispacy: NER,Spacy models for science papers.
- Bokeh: A library for visualising interacted charts.

Data exploration

CORD-19 dataset











-  Kaggle
- ▼  cord_19_embeddings
 -  cord_19_embeddings_...
- ▼  document_parsers
 -  pdf_json
 -  pmc_json
 -  COVID.DATA.LIC.AGMT....
 -  json_schema.txt
 -  metadata.csv
 -  [metadata.readme](#)

Figure 1: CORD-19 Structure

2022-05-12

Text mining for exploration of COVID-19 severity

factors

└ Data exploration

└ Data exploration

Data exploration

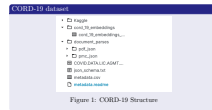


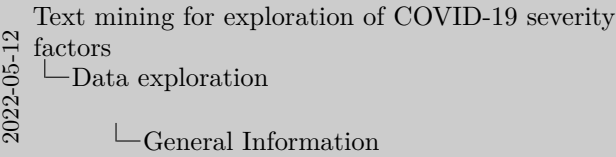
Figure 1: CORD-19 Structure

General Information

The metadata consist of more than one millions articles.

	cord_uid	sha	source_x	title	doi	pmcid	pubmed_id	license	abstract	publish_time	author
0	ug7v899j	d1aafb70c066a2068b02786f8929fd9c900897fb	PMC	Clinical features of culture-proven Mycoplasma...	10.1186/1471-2334-1-6	PMC35282	11472636	no-cc	OBJECTIVE: This retrospective chart review des...	2001-07-04	Madani, Tariq A; Al-Ghamdi Aisha ,
1	02tnwd4m	6b0567729c2143a66d737eb0a2f63f2dce2e5a7d	PMC	Nitric oxide: a pro-inflammatory mediator in l...	10.1186/rr14	PMC59543	11667967	no-cc	Inflammatory diseases of the respiratory tract...	2000-08-15	Vliet, Alber van der Eiserich, Jasco P; Cros.
2	ejv2xln0	06ced00a5fc04215949aa72528f2eeaae1d58927	PMC	Surfactant protein-D and pulmonary host defense	10.1186/rr19	PMC59549	11667972	no-cc	Surfactant protein-D (SP-D) participates in th...	2000-08-25	Crouch, Erik t
3	2b73a28n	348055649b6b8cf2b9a376498df9bf41f7123605	PMC	Role of endothelin-1 in lung disease	10.1186/rr44	PMC59574	11686871	no-cc	Endothelin-1 (ET-1) is a 21 amino acid peptide...	2001-02-22	Fagan, Kare A; McMurtry Ivan F Rodman, David h
4	9785vg6d	5f48f792a5fa08bed9f56016f4981ae2ca6031b32	PMC	Gene expression in epithelial cells in respons...	10.1186/rr61	PMC59580	11686888	no-cc	Respiratory syncytial virus (RSV) and pneumoni...	2001-05-11	Domachowski Joseph E Bonville Cynthia A Ro.

Figure 2: Overview of the meta data of the dataset



General Information

Columns in the metadata

```
['cord_uid', 'sha', 'source_x', 'title', 'doi',
 'pmcid', 'pubmed_id', 'license', 'abstract',
 'publish_time', 'authors', 'journal',
 'mag_id', 'who_covidence_id', 'arxiv_id',
 'pdf_json_files', 'pmc_json_files', 'url', 's2_id']
```

2022-05-12

Text mining for exploration of COVID-19 severity factors

└─Data exploration

└─General Information

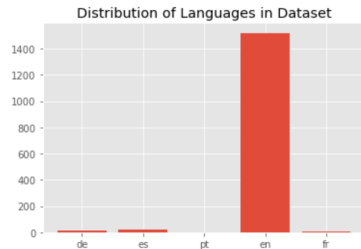
General Information

```
Columns in the metadata
['cord_uid', 'sha', 'source_x', 'title', 'doi',
 'pmcid', 'pubmed_id', 'license', 'abstract',
 'publish_time', 'authors', 'journal',
 'mag_id', 'who_covidence_id', 'arxiv_id',
 'pdf_json_files', 'pmc_json_files', 'url', 's2_id']
```

Language status

As can be observed on a graph most of the papers are written in english. However, there were some exceptions.

During this part, we deleted all non-english articles by using langdetect library.



2022-05-12

Text mining for exploration of COVID-19 severity factors
 └ Data exploration
 └ Language status

Language status

As can be observed on a graph most of the papers are written in english. However, there were some exceptions. During this part, we deleted all non-english articles by using langdetect library.



Data preprocessing

- Converting JSON format into DataFrame format.
- Removing all non-english paper.
- Removing special characters
- Removing numbers
- Tokenizing.
- Removing stopwords.
- Stemming.
- Lemmatisation.

2022-05-12

Text mining for exploration of COVID-19 severity factors

- └ Data preprocessing
 - └ Data preprocessing

Data preprocessing

- Converting JSON format into DataFrame format.
- Removing all non-english paper.
- Removing special characters
- Removing numbers
- Tokenizing.
- Removing stopwords.
- Stemming.
- Lemmatisation.

Data preprocessing

2022-05-12

Text mining for exploration of COVID-19 severity factors

└─Data preprocessing

└─Data preprocessing

Data preprocessing

Data processing

- Data selection
 - Selecting articles with risk factors and severity key-words.
 - Clustering using Latent Dirichlet Allocation.
- NER (Named-entity recognition).

2022-05-12

Text mining for exploration of COVID-19 severity factors

- └ Data processing
 - └ Data processing

Data processing

- Data selection
 - Selecting articles with risk factors and severity key-words.
 - Clustering using Latent Dirichlet Allocation.
- NER (Named-entity recognition).

Introduction	State of the art	Data exploration	Data preprocessing	Data processing	Result
○	○	○○○○	○○	○●○○○○○	○○○○
Risk factors and severity paper filtering					
Dictionary of key words					
Khang Duy LAI - Mariia KLIMINA			Université Paris Cité		
Text mining for exploration of COVID-19 severity factors			12 of 21		

2022-05-12

Text mining for exploration of COVID-19 severity factors

└─Data processing

└─Risk factors and severity paper filtering

Risk factors and severity paper filtering

Dictionary of key words

2022-05-12

- └ Risk factors and severity paper filtering



LDA

In natural language processing, the latent Dirichlet allocation (LDA) is a generative statistical model that allows sets of observations to be explained by unobserved groups that explain why some parts of the data are similar.

The LDA algorithm structure:

- Providing to an algorithm a certain number of topics.
- The algorithm is assigning every word to a temporary topic.
- The algorithm is checking and updating topic assignments.

2022-05-12

Text mining for exploration of COVID-19 severity factors
 └ Data processing
 └ LDA

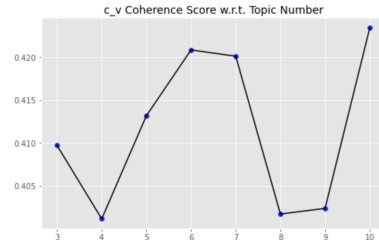
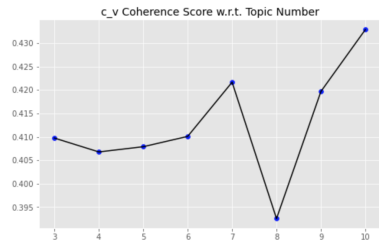
LDA

In natural language processing, the latent Dirichlet allocation (LDA) is a generative statistical model that allows sets of observations to be explained by unobserved groups that explain why some parts of the data are similar.

The LDA algorithm structure:

- Providing to an algorithm a certain number of topics.
- The algorithm is assigning every word to a temporary topic.
- The algorithm is checking and updating topic assignments.

LDA



The coherence score measures how similar these words are to each other. The higher the coherence score is, the more suitable the topic number should be.

2022-05-12

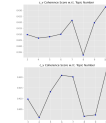
Text mining for exploration of COVID-19 severity

factors

└ Data processing

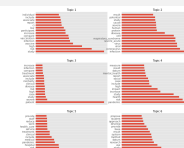
└ LDA

LDA



The coherence score measures how similar these words are to each other. The higher the coherence score is, the more suitable the topic number should be.

2022-05-12

 \perp_{LDA} 

17 of 21

Result

Example of table of result

0	chronic obstructive pulmonary disease copd	DISEASE
1	death	DISEASE
3	copd	DISEASE
9	dyspnea	DISEASE
10	cough	DISEASE
11	copd pulmonary function	DISEASE
13	respiratory tract infection	DISEASE
14	chronic unstable disease system malignancy	DISEASE
19	obstructive pulmonary disease	DISEASE
21	copd airflow	DISEASE
25	hypertension	DISEASE
26	atherosclerotic heart disease	DISEASE

2022-05-12

Text mining for exploration of COVID-19 severity factors

└─Result

└─Result

Result

Example of table of result	
0	chronic obstructive pulmonary disease copd
1	death
3	copd
9	dyspnea
10	cough
11	copd pulmonary function
13	respiratory tract infection
14	chronic unstable disease system malignancy
19	obstructive pulmonary disease
21	copd airflow
25	hypertension
26	atherosclerotic heart disease
27	bronchiectasis

Conclusion

- Basic diseases are filtered out
- The method have not sorted if the disease is a covid symtoms or related to the progress of severity cases.

2022-05-12

Text mining for exploration of COVID-19 severity factors
└ Result
└ Conclusion

Conclusion

- Basic diseases are filtered out
- The method have not sorted if the disease is a covid symtoms or related to the progress of severity cases.

Future improvement

- Create a knowledge graph
- Calculate the severity rate

2022-05-12

Text mining for exploration of COVID-19 severity factors
└ Result
└ Future improvement

Future improvement

- Create a knowledge graph
- Calculate the severity rate

Introduction

State of the art

Data exploration

Data preprocessing

Data processing

Result

○

○

○○○○

○○

○○○○○○○

○○○●

Thank you

Thank you for your listening.

Khang Duy LAI - Mariia KLIMINA

Université Paris Cité

Text mining for exploration of COVID-19 severity factors21 of 21

2022-05-12

Text mining for exploration of COVID-19 severity factors

└─Result

└─Thank you

Thank you

Thank you for your listening.