



# Dictionary-based text tagging of articles related to COVID-19

Jesper Laurell & Jennie Karlsson



# Project idea

- Covid-19 Pandemic
- Kaggle Task
- Open research dataset

Massive  
amount of  
Articles

Text mining

Get  
Overview



# Tagger

## Dictionaries

- Virus
- Disease
- Symptoms

## Input

- Json
- Parsing
- Kaggle database

```
1 covid19
2 covid 19
3 ncp
4 coronavirus disease 2019
5 corona virus disease 2019
6 coronavirus disease19
```

# Tagger

## Output - Json

```
{"cord_uid":"PMC7159299","sourcedb": "PMC", "sourceid": "n/a",  
"div_id": 0, "text": "Clinical features of patients infected  
with 2019 novel coronavirus in Wuhan, China. ", "denotations":  
[{"id": "JJ_dict_v1_virus", "span": {"begin": 44, "end": 66},  
"obj": "Virus_SARS-CoV-2"}]}
```

## PubAnnotation

Clinical features of patients infected with 2019 novel coronavirus in Wuhan, China.

Virus\_SARS-CoV-2



# Tagger - the algorithm

## Idea

- Sort dictionaries
- Regex - finditer
- Loop article for each dict

## Considerations

- Hyphens
- Capital Letters



# Evaluation

- Gold standard
- Compare
- Performance assessment crucial
- Identify weaknesses
- Risk: optimize tagger for gold std



## Evaluation scores

- Recall
- Precision
- Micro
- Macro
- F1 score
- Runtime

$$\text{Recall} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$$

$$\text{Precision} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}}$$

$$F1 = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$



# Results

<b>Dictionaries</b>	<b>Precision</b>	<b>Recall</b>	
Disease_COVID-19	70%	70%	
Virus_SARS_CoV-2	40%	15%	
Symptom_COVID-19	68%	61%	
<b>Total</b>	<b>Precision</b>	<b>Recall</b>	<b>F1-score</b>
Micro	67%	59%	63%
Macro	59%	49%	53%



# Results

Class	TP	FP	FN	True Total
Disease_COVID-19	32	14	14	46
Virus_SARS-CoV-2	2	3	11	13
Symptom_COVID-19	17	8	11	28

	Runtime [s]
Tagging gold std	1.5
Evaluating gold std	0.002
Tagging subset_100	120



## Final version

### Assessment

- Runtime
- Compact code
- 2/3 dicts acceptable results

### Future improvements

- Improve dictionaries
- Use more dictionaries



**THANK YOU!**

**QUESTIONS?**