# *From Bats to Masks:*
# Change of Topics in Swedish Articles about Coronavirus

Bernadeta Griciute[1,2], Lifeng Han[3], Goran Nenadic[3]

University of Malta[1], Saarland University[2], University of Manchester[3]

# Content

Motivation

Corpus preparation

Experimental work using LDA tools

Findings

# Topic Modelling of Swedish Newspaper Articles about Coronavirus: a Case Study using Latent Dirichlet Allocation Method

Bernadeta Griciūtė, Lifeng Han, Goran Nenadic

Topic Modelling (TM) is from the research branches of natural language understanding (NLU) and natural language processing (NLP) that is to facilitate insightful analysis from large documents and datasets, such as a summarisation of main topics and the topic changes. This kind of discovery is getting more popular in real-life applications due to its impact on big data analytics. In this study, from the social-media and healthcare domain, we apply popular Latent Dirichlet Allocation (LDA) methods to model the topic changes in Swedish newspaper articles about Coronavirus. We describe the corpus we created including 6515 articles, methods applied, and statistics on topic changes over approximately 1 year and two months period of time from 17th January 2020 to 13th March 2021. We hope this work can be an asset for grounding applications of topic modelling and can be inspiring for similar case studies in an era with pandemics, to support socio-economic impact research as well as clinical and healthcare analytics. Our data and source code are openly available at https://github. com/poethan/Swed_Covid_TM Keywords: Latent Dirichlet Allocation (LDA); Topic Modelling; Coronavirus; Pandemics; Natural Language Understanding

https://doi.org/10.48550/arXiv.2301.03029

# Motivation

## China Grapples With Mystery Pneumonia-Like Illness

Beijing is racing to identify a new illness that has sickened 59 people as it tries to calm a nervous public.
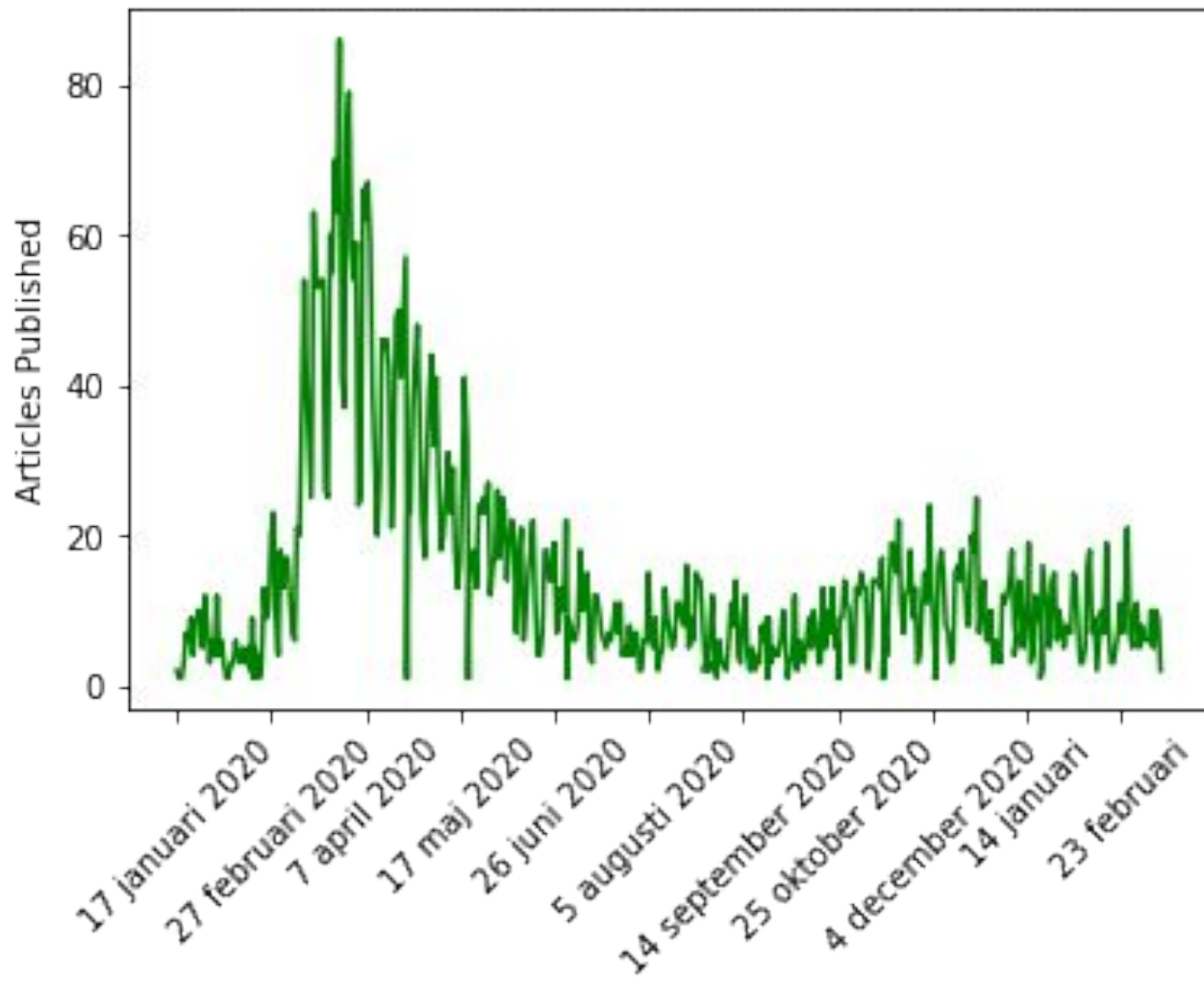
Published Jan. 6, 2020

## US declares public health emergency over coronavirus, announces temporary travel ban

The seventh U.S. coronavirus case has been reported in California.

February 1, 2020, 6:40 PM
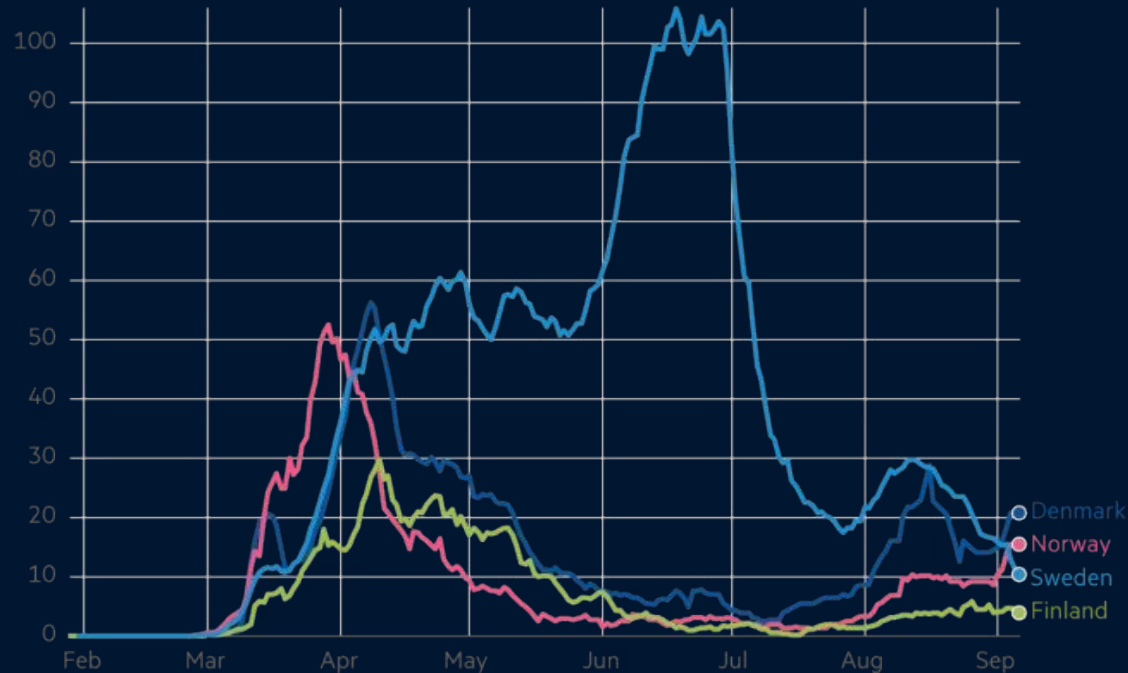
# How many Corona-related articles per day?

Articles Published

17 januari 2020
27 februari 2020
7 april 2020
17 maj 2020
26 juni 2020
5 augusti 2020
14 september 2020
25 oktober 2020
4 december 2020
14 januari
23 februari

SVT, Sweden

# Different approaches

The profile of Sweden's pandemic differs radically from those of its neighbours

New confirmed cases of Covid-19 (per million)

New confirmed cases of Covid-19, seven-day rolling average of new cases (per million)

# Topic Modelling.
# Latent Dirichlet Allocation (LDA)

# More LDA

$$p(\beta_{1:K}, \theta_{1:D}, z_{1:D}, w_{1:D})$$
$$= \Pi_{i=1}^{K} p(\beta_i) \Pi_{d=1}^{D} p(\theta_d)$$
$$\left( \Pi_{n=1}^{N} p(z_{d,n}|\theta_d) p(w_{d,n}|\beta_{1:K}, z_{d,n}) \right)$$

where the four main parameters $\beta$, $\theta$, $z$, and $w$ represent respectively the "topic distribution", "topic proportion of document", "topic assignment of document", and the "observed words of document".

(Blei et al., 2003; Blei, 2012)

Read our paper for more detailed interpretations https://arxiv.org/abs/2301.03029

# DTM

In comparison to "statistical assumptions of a **static** topic model, such as Latent Dirichlet Allocation (LDA) (Blei et al., 2003)."

- LDA assumes the documents are drawn exchangeably form the same set of topics.
- However, the order of some collections reflects an evolving set of topics

To address this, DTM approaches the task by dividing data by time slice, e.g. day/month/year

- Model the documents with k-component topic model
- Topics associated with slit t evolve forms the topics associated with slice 't-1'
- Read Blei and Lafferty (2006) for more math implementation, included in **Gensim toolkit**

Dynamic Topic Models (DTM), introduced by Blei and Lafferty (2006); Blei (2012)

# Gensim



Gensim is a FREE Python library

# Topic modelling for humans

✓ Train large-scale semantic NLP models

✓ Represent text as semantic vectors

✓ Find semantically related documents

# Choosing Number of Topics

# Dataset

Newspaper articles having Covid as one of the main topics

SVT (Sveriges Television) - Sweden's national public television broadcaster

2 251 article

One year from the first article: 2020/01/17 - 2021/01/17

# Dataset: how to get it?



https://github.com/poethan/Swed_Covid_TM

# Findings

# From Local to Global Issue



Coronavirus
China
Chinese
country
person
people
infect
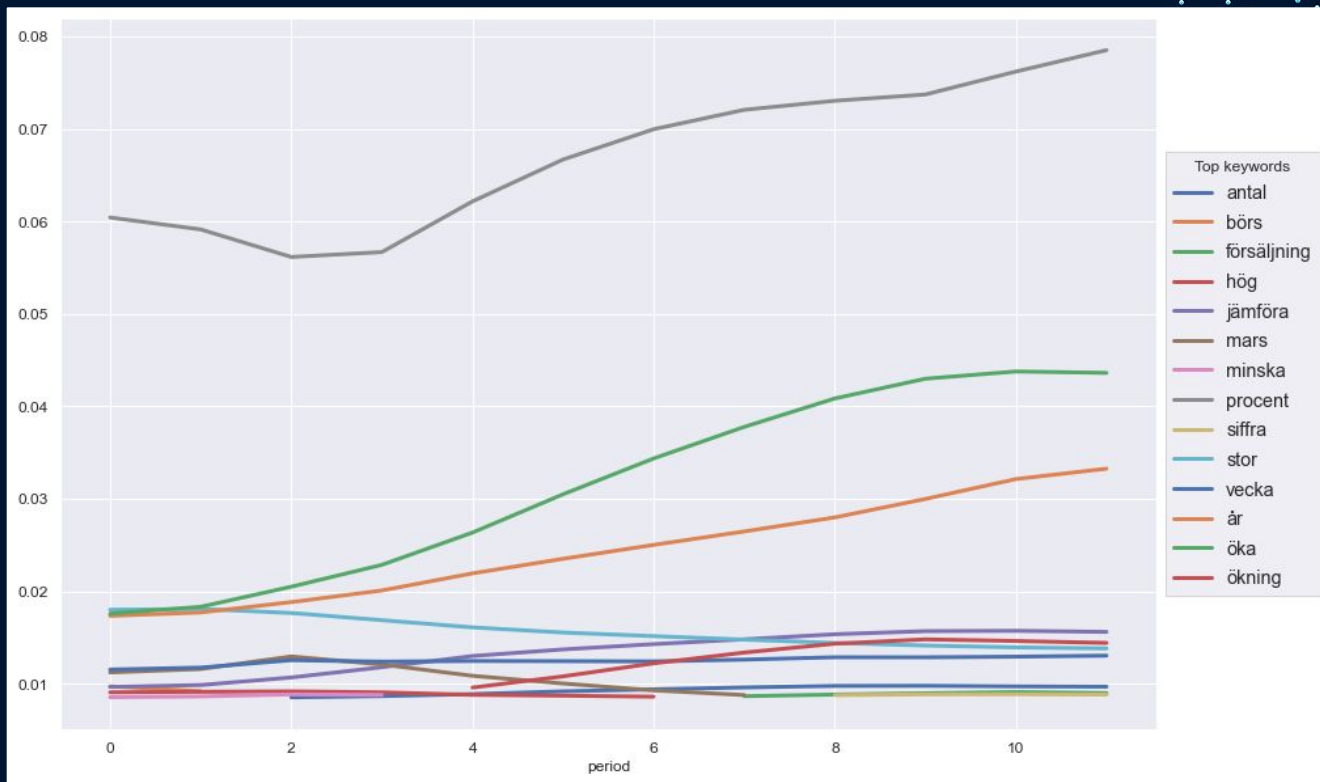spread
outbreak
virus
world
WHO
Wuhan

# Factual Information

# From Recommendations to Advices



general
distance
Public Health Agency
apply to
decrease
face mask
authority
person
recommendation
travel
risk
advice
spread of infection
big
avoid

# Consequences for Economy



number
stock market
sales
high
compare
March
decrease
percent
numbers
big
week
year
increase

# Data with code:



https://github.com/poethan/Swed_Covid_TM

# Future Work

Expand the data frame of the data

Use more varied data sources

Compare to other languages

# Thank you!

Time for questions and suggestions :)