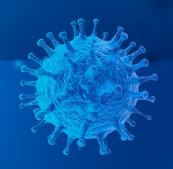


Covid-19 News Data

Team Torreys Peak **Dylan Bernstein, Jennifer Dickson, Katie Greenfield, Madison Moye, and Yongbo Shu**



The Data

COVID-19 News Articles

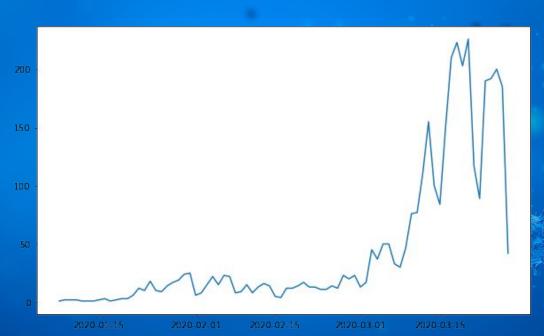
- → 3500+ CBC News Articles
 - CBC = Canadian Broadcasting Company
- Contains authors, title, publish date, description, full article, and URL.
- → Articles date from Dec, 2019 Mar 27, 2020

https://www.kaggle.com/ryanxjhan/cbc-news-coronavirus-articles-march-26

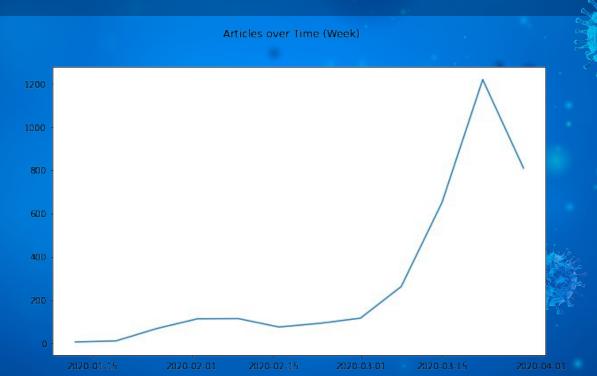


The Data





The Data



Three Approaches

Word Cloud

Data Visualization tool that summarizes key topics or words in data in a visually intuitive manner.

Topic Modeling

Unsupervised machine learning model that deciphers topics in a corpus of data or documents.

Word Count

Using a pySpark RDD and pySpark Dataframe to count all the occurrences of words used in the articles.

1. Word Cloud

Word Cloud

Created word clouds for all text in the articles and also created word clouds for December, January, February and March to see if the articles had changed over time.

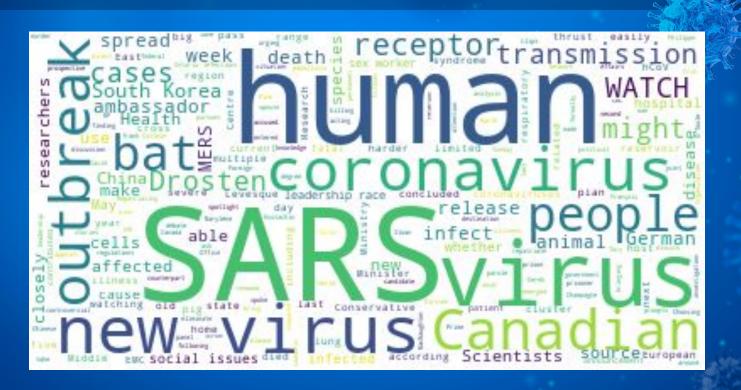
Found 15.8M words in the text of the articles.



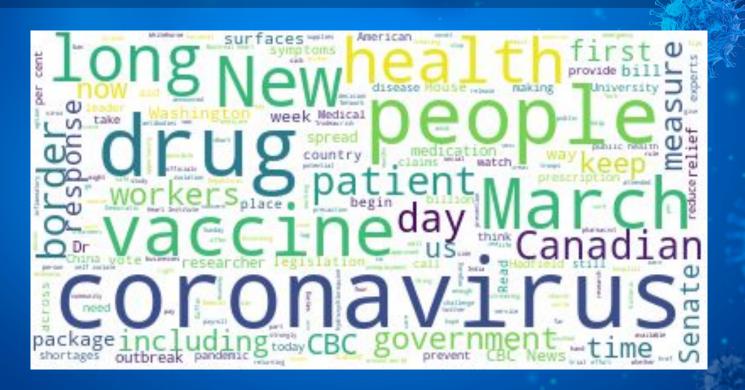
Word Cloud - All Articles



Word Cloud - December 2019



Word Cloud - March 2020



2. Topic Modeling

Topic Modeling — Data Preprocessing

Tokenizing

"Bob walked the dog" →
[Bob, walked, the, dog]

Remove Special Characters

!,";-?- etc.

Position Tagging

Noun, verb, adj

Remove Stopwords

'http' or 'nt'

Lemmatizing

"Walked" → "walk"

Two Topic Models

"Bigger" – less data preprocessing

"Smaller" – much more data preprocessing

"Bigger" Topic Model

Topic 1:	Topic 2:	Topic 3:	Topic 4:	Topic 5:	Topic 6:
#1: pei	#1:	#1:	#1:	#1: per	#1: flight
#2: islanders	#2: the	#2: the	#2: the	#2: cent	#2: travel
#3: morrison	#3: be	#3: to	#3: be	#3: price	#3: airline
Topic 7:	Topic 8:	Topic 9:	Topic 10:	Topic 11:	Topic 12:
#1:	#1: community	#1: http	#1: mask	#1: sanders	#1: party
#2: the	#2: territory	#2: href	#2: hand	#2: biden	#2: budget
#3: be	#3: nwt	#3: 5ewtf	#3: supply	#3: read	#3: election

"Smaller" Preprocessing

- Remove numbers and short words
- Remove common tokens (Threshold = 0.9)
- Remove Uncommon Tokens (Threshold=0.01)

▶ Beta of 0.5



"Smaller" Topic Model

Topic 1:	Topic 2:	Topic 3:	Topic 4:	Topic 5:	Topic 6:
#1: people	#1: bc	#1: pei	#1: ship	#1: government	#1: health
#2: masks	#2: people	#2: nova	#2: cruise	#2: trump	#2: case
#3: health	#3: health	#3: scotia	#3: passenger	#3: party	#3: hospital
Topic 7:	Topic 8:	Topic 9:	Topic 10:	Topic 11:	Topic 12:
#1: people	#1: canada	#1: case	#1: school	#1: health	#1: business
#2: city	#2: flight	#2: virus	#2: student	#2: case	#2: cent
#3: service	#3: canadians	#3: china	#3: time	#3: province	#3: government

Counted occurrences of all words that were used in all these news articles.

Two ways:

- pySpark rdd
- pySpark dataframe.



complying with the Canadian Pharmacists Association call to limit the amount of medications given to patients to 30-d ay\xa0supplies. The goal is to stop people from refilling prescriptions early and to ensure life-saving drugs don\'t run short when supply chains are\xa0vulnerable. Mina Tadrous is a pharmacist and researcher in Toronto who monitors p

```
newdf = df.withColumn('text', regexp_replace('text', '\xa0', ''))
newdf = newdf.withColumn('text', regexp_replace('text', '\\\'', ''))
```

```
[('\U0001f9d0)', 1),
 ('\U0001f970\\', 1),
 (' @ — @RubyTuesday 72Tried', 1),
 ('\efticolemnic <a', 1),
 (' & & mdash; @hannahmaryr', 1),
 ('\equiv ', 1),
 (' * — @jdutchermusic', 2),
 ('\verts', 1),
 ('\varphi', 3),
 ('  — @DocSchmadia', 2),
 (' 👋 ', 3),
 (' 4 4 ', 1),
 (' - < a', 1),
 (' - <br>We', 1),
 ('^{*})', 1),
 (' | <br>>Meet', 1),
 ('\ufeffCorriveau\ufeff', 1),
 ('-', 2),
 ('\begin{aligned} <a', 1)]
```

- --- rdd running time (without sorting) is 0.030278921127319336 seconds ----- DataFrame running time (without sorting) is 0.18393468856811523 seconds ----- rdd running time (with sorting) is 8.592974662780762 seconds ---
- --- DataFrame running time (with sorting) is 0.2873249053955078 seconds ---
- --- DataFrame running time (with sorting) is 0.28/32490539550/8 seconds ---

- --- rdd running time is 0.7412478923797607 seconds ---
- --- DataFrame running time is 8.040996551513672 seconds ---
- --- DataFrame running time is 1.1875431537628174 seconds ---

['Cbc News'] Coronavirus a 'wa 2020-03-27 08:00:00 Canadian pharmaci https://www.cbc.c 36
['The Associated U.S. Senate passe 2020-03-26 05:13:00 The Senate late W https://www.cbc.c 9' ['Cbc News'] Coronavirus: The 2020-03-27 00:36:00 Scientists around https://www.cbc.c 7' ['Cbc News'] The latest on the 2020-03-26 20:57:00 Trudeau says r https://www.cbc.c 13' ['Mark Gollom Is 'Worse' pandemic 2020-03-27 08:00:00 The continued exi https://www.cbc.c 10' ['Cbc News'] What you need to 2020-03-27 08:00:00 Recent developmen https://www.cbc.c 18'
['Cbc News'] Coronavirus: The 2020-03-27 00:36:00 Scientists around https://www.cbc.c 73 ['Cbc News'] The latest on the 2020-03-26 20:57:00 Trudeau says r https://www.cbc.c 133 ['Mark Gollom Is 'Worse' pandemic 2020-03-27 08:00:00 The continued exi https://www.cbc.c 103 ['Cbc News'] What you need to 2020-03-27 08:00:00 Recent developmen https://www.cbc.c 183 Recent developmen https://www.cbc.c https://www.cbc.c https://www.cbc.c https://www.cbc.c https://www.cbc.c https://www.cbc.c https://www.cbc.c https:
['Cbc News'] The latest on the 2020-03-26 20:57:00 Trudeau says r https://www.cbc.c 136 ['Mark Gollom Is 'Worse' pandemic 2020-03-27 08:00:00 The continued exi https://www.cbc.c 105 ['Cbc News'] What you need to 2020-03-27 08:00:00 Recent developmen https://www.cbc.c 182 183
['Mark Gollom Is 'Worse' pandemic 2020-03-27 08:00:00 The continued exi https://www.cbc.c 10:00
['Cbc News'] What you need to 2020-03-27 08:00:00 Recent developmen https://www.cbc.c 182
['The Associated Michigan hospital 2020-03-26 11:02:00 Michigan hospital https://www.cbc.c
['Thomson Reuters'] U.S. coronavirus 2020-03-26 14:55:00 The number of con https://www.cbc.c 14
['Leah Hendry Is 'Avoid the emerge 2020-03-27 08:00:00 The Jewish Genera https://www.cbc.c 6
['Reporter', 'Web COVID-19 in Sask: 2020-03-26 14:18:00 Three Saskatchew https://www.cbc.c
['Jorge Barrera I Manitoba chiefs o 2020-03-27 08:01:00 A Manitoba chiefs https://www.cbc.c
['Colleen M. Floo How invoking the 2020-03-27 08:00:00 This column is an https://www.cbc.c 10
['Producer', 'Cbc In Ontario, const 2020-03-27 08:00:00 Construction is o https://www.cbc.c
['Dan Mcgarvey Is Alberta's film in 2020-03-26 12:00:00 Albertas TV and m https://www.cbc.c 68
['Cbc News New Yo 'Like a war zone' 2020-03-27 08:00:00 The first wave hi https://www.cbc.c
[] N.L. fisheries re 2020-03-27 08:30:00 A Memorial Univer https://www.cbc.c 49
[] 1st death, 3 new 2020-03-11 00:15:00 Manitoba is under https://www.cbc.c 478
['Investigative R Medical experts w 2020-03-26 08:00:00 Medical experts a https://www.cbc.c
['Hadeel Ibrahim It's 'too late' f 2020-03-27 08:00:00 New Brunswick has https://www.cbc.c 123
+++++++

-20

Ingest data to S3 bucket

```
counts_rdd.write.parquet('s3://msbx5420-2020/Team-Torreys-Peak/counts.parquet')
nnewdf.write.parquet('s3://msbx5420-2020/Team-Torreys-Peak/news_countnews.parquet')
[[hadoop@ip-172-16-1-40 ~]$ aws s3 ls s3://msbx5420-2020/Team-Torreys-Peak/PRE counts.parquet/
```

PRE news_countnews.parquet/

2020-04-26 04:38:51 17478084 news.csv

Thank you!

Questions?

