

## Cyberpsychology, Behavior, and Social Networking

Cyberpsychology, Behavior, and Social Networking: <http://mc.manuscriptcentral.com/cyberpsych>

### **#coronavirus: Monitoring the Belgian Twitter Discourse on the Sars-Cov-2 Pandemic**

Journal:	<i>Cyberpsychology, Behavior, and Social Networking</i>
Manuscript ID	CYBER-2020-0341.R1
Manuscript Type:	Original Article
Keyword:	Twitter, Quantitative Research
Manuscript Keywords (Search Terms):	Sars-coV-2, pandemic, covid-19, Twitter, Belgium

SCHOLARONE™  
Manuscripts

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

**#coronavirus:**

**Monitoring the Belgian Twitter Discourse on the SARS-CoV-2 Pandemic**

**Running Title: Twitter Discourse on SARS-CoV-2**

For Peer Review ONLY/Not for Distribution

**Running head: TWITTER DISCOURSE ON SARS-COV-2***Abstract*

In the current study, a social media analysis is conducted to examine the public discourse about the SARS-CoV-2 pandemic on Twitter. In particular, this study aims to examine: (1) How the number of tweets varies as a function of the timeline of the pandemic and associated measures, (2) How the content of these tweets, including displayed emotions, changes. Therefore, 373,908 tweets and retweets from Belgium were collected from the 25<sup>th</sup> of February 2020 to the 30<sup>th</sup> of March. Time series analysis, network bigrams, topic models, and emotional lexica were deployed for analysis. The results showed that significant events related to the virus correlated with an immediate increase in the number of tweets addressing them. Furthermore, the Belgian Twitter discourse was characterized by positively connotated words, which also refer to European solidarity. These findings do not only stress the relevance of Twitter as a medium for public discourse during lockdowns, but also seem to indicate that the Belgian public supports policy measures which respect solidarity in Europe.

*Keywords:* SARS-CoV-2, pandemic, covid-19, Twitter, Belgium

**Running head: TWITTER DISCOURSE ON SARS-COV-2****Introduction**

The recent outbreak of the SARS-CoV-2 pandemic has disrupted the public life across the world in a manner unseen since the end of the Second World War. The virus, first identified in Wuhan (China) in December 2019, was recognized by the World Health Organization as a pandemic in March 2020 causing a severe threat to the health of people around the world. Given the serious public health risk, it was of critical importance for policy makers to respond early to be able to slow down the community spread of the virus and as such lower the burden on healthcare system. Following the examples of other countries, the Belgian government decided in several stages to take far-reaching measures to curtail public life and urge people to stay at home to further decrease the spread of the virus. In the current study, a social media analysis is conducted to examine the public discourse about the SARS-CoV-2 pandemic on Twitter. In particular, this study wants to answer two research questions: (1) How does the number of tweets vary as a function of the timeline of the pandemic and associated measures, (2) How does the content of these tweets, including displayed emotions, vary over time. Gaining profound insight in how the public perceives this pandemic and how it reacts to the measures declared by the government is of utmost importance for policymakers, health workers, and stakeholders who communicate to the public during infectious disease outbreaks.

Agenda setting theory describes the ability of media to alter the salience of issues on the public agenda<sup>1</sup>. According to basic agenda setting, issues that get a lot of attention in the media are considered more important by the public<sup>2</sup>. It has been argued that nowadays, Twitter has at least as much influence on the public agenda as newspapers<sup>3,4</sup>. During health crises, Twitter gives users the possibility to express opinions and share breaking news<sup>5</sup>. Initial studies on tweets during emerging epidemics have found that Twitter can be used to identify public concerns<sup>6</sup> or even trace infectious activity<sup>7,8</sup>. Furthermore, it can be an indicator for the public risk perception of an infection<sup>9</sup>. Although these initial studies report on the use of Twitter

**Running head: TWITTER DISCOURSE ON SARS-COV-2**

1  
2  
3 during infectious disease outbreaks like Zika<sup>6,9</sup>, Bird Flu (H7N9)<sup>5</sup> and Swine Flu (H1N1)<sup>7,8</sup>,  
4  
5 examining the public discourse on Twitter during the SARS-CoV-2 pandemic is particularly  
6  
7 relevant given that the scale of it was unseen and it had an enormous impact on the daily lives  
8  
9 of billions of people around the world.  
10  
11  
12

**Method****Sample**

13  
14  
15  
16  
17  
18 For the current study, corona-related tweets were collected from February 25, 2020 to  
19  
20 March 30. The first patient in Belgium tested positive on SARS-CoV-2 on the 4<sup>th</sup> of February;  
21  
22 however press attention only started to rise after about 100 Belgians were quarantined in an  
23  
24 hotel on Tenerife, from the 25<sup>th</sup> of February onwards<sup>10,11</sup>. This event marks the start of our data  
25  
26 collection. Thus, the dataset in the current study provides an overview of corona-related tweets  
27  
28 during the unfolding of the SARS-CoV-2 pandemic in Belgium.  
29  
30  
31

32 At the start of the data collection, it appeared that the hashtags ‘coronavirus’, ‘covid19’,  
33  
34 ‘coronavid19’, and ‘corona’ were used most frequently in Belgium. Resultantly, these hashtags  
35  
36 were used to collect the data. However, starting on March 11 ‘Coronavirusoutbreak’,  
37  
38 ‘COVID19BE’, and ‘coronabelgie’ were also regarded because Twitter displayed them as  
39  
40 trending. Tweets were collected through the streaming Twitter API access in combination with  
41  
42 the twitterR<sup>12</sup> package in R.  
43  
44  
45

46 All tweets and retweets posted in Belgium in Dutch, French, and English language were  
47  
48 included in the data collection, although the streaming Twitter API only allowed us to capture  
49  
50 the first 128 characters of a sample of actual posted tweets. In total, the scraping procedure  
51  
52 captured 373.908 tweets and retweets of which 71.481 were original tweets. A subsample of  
53  
54 31.454 were written in English, 15.936 in French, and 24.091 in Dutch.  
55  
56  
57  
58  
59  
60

**Running head: TWITTER DISCOURSE ON SARS-COV-2****Analysis**

Prior to analysis, tweets were cleaned and stemmed using the text mining<sup>13</sup> and tidytext<sup>14</sup> packages. In particular, stop words, punctuation and other special signs like smileys, hashtags, slashes etc. were removed. Porter's Snowball stemmer<sup>15</sup> was applied to reduce the processed words to their word stems (e.g., reducing 'washing' to 'wash'). In order to answer RQ1, the number of tweets and retweets were computed on an hourly basis for the given timeframe. To tackle RQ2, three analytical techniques were applied. First, network bigrams examined which pairs of words frequently appear together and how they relate to each other. Second, to analyze the emotional connotation of the words, an emotion lexicon based on mechanical turk<sup>16</sup> was used. Third, machine learning based topic models were constructed to show the themes which the tweets addressed frequently. Latent dirichlet allocation was used for that purpose as it is implemented in the topic models package<sup>17</sup> in R.

**Results****Evolution of tweets over time**

Figure 1 contains the number of English, Dutch, and French tweets and retweets on the given hashtags aggregated per hour as received by the streaming API. The first peak occurs on March 2. On this day, the number of infections jumped from two to eight in Belgium<sup>18</sup>. Afterwards the number of tweets decreased to its initial level again until another spike occurred on Monday March 9 when stock markets crashed worldwide leading to a loss of the BEL20 of about 8%. The first death resulting from SARS-CoV-2 in Belgium was announced on the 11<sup>th</sup> of March by the Belgian Minister of Health<sup>19</sup>. The Government announced its decision to close schools on March 12<sup>20</sup>, which correlated with an increase in the number of tweets. It was communicated that the schools should stay closed at least until the end of the Easter holidays on April 20. On the 14th and 15th of March the number of tweets dropped again. This is because the number of tweets is lower on the weekends. This trend of lower tweets on weekends

**Running head: TWITTER DISCOURSE ON SARS-COV-2**

continues throughout the time series. On the 17<sup>th</sup> of March, stricter social distancing measures were introduced<sup>21</sup>, meaning that leaving the place of residence was restricted to essential activities like grocery shopping, medical treatments, and exercise. Furthermore, it was generally restricted to meet in groups larger than two, except for direct relatives living in the same household. The spike on March 17 corresponds to the hour when the Belgian prime minister publicly announced the implementation of that measure. The highest number of tweets could be registered on March 19, when Michel Barnier, a well-known French politician, announced via Twitter that he had tested positive for the virus. Another rise in the number of tweets occurred one week later when the US congress approved a two trillion dollar economical stimulus<sup>22</sup>. On March 30, a high number of tweets addressed that the Hungarian government passed an emergency bill, resulting in a variety of executive rights for the president.

It can also be stated that the number of tweets increased rapidly after the mentioned happenings took place. Most of the tweets were sent within two hours after the occasion. The exception was the US approval of the 2 trillion dollar stimulus which became public between 5 and 6 AM Belgian time which was associated with a later spike of tweets between 10 and 11 AM.

[FIGURE 1 HERE]

**Content of the tweets**

For further semantic analysis, we removed the retweets from the dataset and focused on tweets in English language. We did so, because this was the most frequent language used within the captured tweets. To answer RQ2 and gain insight in the content of the tweets, a network of bigrams was computed. The network shows, which terms frequently occur together. By applying that technique, we can learn how different words are combined within tweets. As shown in Figure 2, some bigrams appear to be largely unrelated to other word pairs, e.g. ‘toilet

**Running head: TWITTER DISCOURSE ON SARS-COV-2**

1  
2  
3 paper' and 'South Korea'. The main network of bigrams is centered around the words 'covid19'  
4  
5 and 'coronavirus', containing words like 'fight', 'spread', 'response', 'test', and 'lockdown'.  
6  
7 The network indicates some positive themes within tweets. For instance, the word 'Italian'  
8  
9 occurs frequently together with 'support' and 'friend'. Also 'stay' and 'inside' are connected  
10  
11 to 'safe'. Although, the word 'death' occurs in the network, there were no agonizing words like  
12  
13 'fear', 'disaster', or similar frequently used.  
14  
15  
16  
17  
18

19 [FIGURE 2 HERE]  
20  
21  
22  
23

24 To gain a further understanding of word relations within the tweets, topic modelling was  
25  
26 deployed. Its aim was to reveal several latent topics and the likelihood of words to appear  
27  
28 within these topics. To select the appropriate number of topics, we followed the guidelines of  
29  
30 Maier et al.<sup>23</sup> and not only used different metrics<sup>24,25,26</sup> but also tenfold cross validation. The  
31  
32 indices showed that at least six different topics should be modelled.  
33  
34  
35  
36  
37

38 [FIGURE 3 HERE]  
39  
40  
41  
42  
43

44 Figure 3 displays the word stems on its y-axis and the density of terms within that topic  
45  
46 on its x-axis represented by beta. It is not surprising that the terms 'coronavirus' and 'covid19'  
47  
48 have a high likelihood to occur in each of the six topics. However, all of these topics stress  
49  
50 slightly different aspects of the discourse. For instance, a tweet addressing the first topic is  
51  
52 likely to contain the term 'crisis' ( $\beta_{\text{crisis1}} = .0093$ ), 'country' ( $\beta_{\text{countri1}} = .0066$ ), 'pandemic'  
53  
54 ( $\beta_{\text{pandem1}} = .0059$ ) and stresses the role of the people to help ( $\beta_{\text{help1}} = .0062$ ) each other and fight  
55  
56 ( $\beta_{\text{fight1}} = .0052$ ) the pandemic. The second topic is also likely to contain the term 'people' ( $\beta_{\text{peopl2}}$ )  
57  
58  
59  
60



## Running head: TWITTER DISCOURSE ON SARS-COV-2

= .0061), but further stresses the need ( $\beta_{\text{need}2} = .0062$ ) for support ( $\beta_{\text{support}2} = .0052$ ). It also uses the word ‘health’ ( $\beta_{\text{health}2} = .0044$ ) and addresses the situation in Italy ( $\beta_{\text{itali}2} = .0036$ ) (which suffered the most from covid19 during the time of the data collection). The fifth topic is similar to the second one. It uses terms like ‘support’ ( $\beta_{\text{support}5} = .0057$ ) and ‘itali’ ( $\beta_{\text{itali}5} = .0041$ ), although it seems to use words with a positive connotation such as ‘good’ ( $\beta_{\text{good}5} = .0041$ ). The sixth and last topic addresses the need ( $\beta_{\text{need}6} = .0068$ ) for Europe wide ( $\beta_{\text{european}6} = .0057$ ) testing ( $\beta_{\text{test}6} = .0067$ ) in fighting ( $\beta_{\text{fight}6} = .0059$ ) the spread ( $\beta_{\text{spread}6} = .0049$ ) of the SARS-CoV-2.

The topic models show the central motives addressed by the tweets. It can be concluded that most of them address the need for European collaboration to fight the spread of the pandemic. Several aspects of containing it are discussed in the tweets like ‘the need for further testing’ in topic 6 and ‘the importance of time’ in topic 3. Furthermore, topics 2 and 5 stress the need for international support for countries who have been severely hit by the pandemic, such as Italy.

### Displayed emotions over time

Finally, the emotions associated with the words within the tweets were examined. They were coded by using an emotional lexicon by Mohammad and Turney<sup>16</sup>. The results are presented in Figure 4. Every layer within the area plot represents the relative amount of a given emotion within original tweets on a daily basis.

Novel to this analysis is the insight which emotions were expressed the most frequently throughout the tweets. On average, the tweets contained terms associated with ‘trust’ the most often ( $\bar{x}_{\text{trust}} = 28.9$  % of all words,  $SD_{\text{trust}} = .03$ ), followed by ‘anticipation’ ( $\bar{x}_{\text{anticipation}} = 17.5$  %,  $SD_{\text{anticipation}} = .02$ ), and ‘fear’ ( $\bar{x}_{\text{fear}} = 17.0$  %,  $SD_{\text{fear}} = .02$ ). Terms which were associated with ‘anger’ ( $\bar{x}_{\text{anger}} = 7.8$  %,  $SD_{\text{anger}} = .01$ ), ‘surprise’ ( $\bar{x}_{\text{surprise}} = 6.4$  %,  $SD_{\text{surprise}} = .01$ ) and ‘disgust’ ( $\bar{x}_{\text{disgust}} = 4.4$  %,  $SD_{\text{disgust}} = .01$ ) occurred the least frequently. In addition, the analysis

**Running head: TWITTER DISCOURSE ON SARS-COV-2**

1  
2  
3 showed that trust had the highest volatility among the emotions. This means that when the  
4  
5 number of tweets increased, the proportion of trustful tweets increased more strongly than  
6  
7 tweets expressing different emotions. The proportion of trustful tweets ranged from 21.5 % on  
8  
9 the 29<sup>th</sup> of February to 37.1 % on the 2<sup>nd</sup> of March. Meaning that the highest proportion of  
10  
11 tweets containing trust was obtained right after the first jump in infections. Trustful tweets  
12  
13 decreased after both the stock market crash on the 9<sup>th</sup> of March and the decision to close schools  
14  
15 on the 12<sup>th</sup> of March. Never-the-less, trustful tweets then rose to their second highest value on  
16  
17 the 16<sup>th</sup> of March (33.8 %). When stricter social distancing measures were introduced on the  
18  
19 following day, the proportion of trustful tweets decreased only slightly on 17<sup>th</sup> (30.7 %) and  
20  
21 18<sup>th</sup> (29.5 %) of March. When the amount of tweets peaked again on the 26<sup>th</sup> a higher  
22  
23 proportion of trust could be detected (32 %), decreasing in the following days.  
24  
25  
26  
27  
28  
29  
30

31 [FIGURE 4 HERE]  
32  
33

**Discussion**

34  
35  
36 The SARS-CoV-2 pandemic was the cause for a curtailment of the public life in  
37  
38 Belgium and many countries worldwide. Governments introduced strict social distancing  
39  
40 measures and prohibited unnecessary displacements. In these situations, social media remains  
41  
42 one of the few areas in which public discourse is still possible. Looking at tweets from the first  
43  
44 month of the unfolding of the pandemic, this study wanted to (1) analyze how the number of  
45  
46 pandemic related tweets changed as a function of critical incidents, and (2) examine the content  
47  
48 of these tweets in terms of its semantic and emotional connotation. This is especially relevant  
49  
50 as it has been shown that Twitter is an important source of information in times of crises<sup>27,28,29</sup>.  
51  
52 As such, this study provides a unique insight in the topics that were discussed by the public  
53  
54 during the start of the SARS-CoV-2 outbreak, in the way the discourse was held, and how  
55  
56 social media users responded to the measures taken by the government.  
57  
58  
59  
60

**Running head: TWITTER DISCOURSE ON SARS-COV-2**

Time series analysis showed that significant events related to the virus resulted in an immediate increase in the number of tweets addressing them. Similar patterns have been observed in past research. Szomszor et al.<sup>30</sup>, for instance, reported substantial increases in the number of tweets in response to a natural hazard, Terpstra et al.<sup>28</sup> made similar observations related to the swine flu pandemic in 2009. The current study showed that the number of tweets peaked in the hours following the governments' declaration of the measurement taken to prevent the further spread of the SARS-CoV-2, including the closing of the schools and stricter social distancing measures. In addition, high numbers of tweets were, amongst others, associated with the worldwide crash of stock markets on the 9<sup>th</sup> of March and Michel Barnier's confirmed infection on the 19<sup>th</sup> of March. These findings further stress the relevance of Twitter as a medium for public discourse during lockdowns and isolations measures; and its potential to influence the salience of the topics on the public agenda. The pace with which the discourse unfolded on Twitter after an event provided some hints for Twitter's role in the agenda setting process. Individuals' appear to turn to Twitter to communicate important events during crises. This effect could have been reinforced by the Belgian lockdown, which made it harder for citizens to discuss relevant issues offline. This implies that examining these tweets provides important information for policy makers, as the tweets help them understand the public response to policy measures.

The content of the tweets was by definition centered on the SARS-CoV-2 pandemic, but stressed different aspects. First, the bigram network and the topic models showed that the Belgian Twitter discourse was characterized by positive connotated words referring to European solidarity. For instance, 'support' and 'European' were prominent terms within the topics, as well as 'friend', 'Italian' and 'support'. These findings seem to indicate that the Belgian public supports policy measures that respect solidarity with European countries suffering from more severe consequences of the pandemic.

**Running head: TWITTER DISCOURSE ON SARS-COV-2**

1  
2  
3 The conclusion that the Belgian Twitter users convey trust in their ability to overcome  
4 the pandemic is further supported by the results from the analysis of emotions within tweets.  
5  
6 Not only was the most frequently expressed emotion ‘trust’, but the relative proportion of trust  
7  
8 increased, when the number of original tweets rose. Although, previous research pointed out  
9  
10 that Twitter can potentially create panic, e.g. during the Ebola epidemic<sup>31</sup>, the results seem to  
11  
12 indicate that the proportion of fear within tweets was rather low. These findings provide some  
13  
14 evidence for Twitter’s attribute agenda setting capabilities. Despite rising case numbers, the  
15  
16 discourse emphasized trustful attributes of the topics on public agenda. We conclude that  
17  
18 positive emotions within the Belgian Twitter discourse have the potential to increase the  
19  
20 importance of positive aspects of the pandemic related agenda. In addition, the trust of the  
21  
22 Belgian public in the ability to cope with the pandemic remained largely intact.  
23  
24  
25  
26  
27

28 Like all research, this study has a number of limitations. First, the streaming Twitter  
29  
30 API only provided a sample of up to 1% of all tweets. Although, it seems likely that the sample  
31  
32 is representative, Twitter does not disclose their sampling procedure and therefore it remains  
33  
34 unknown. Second, while English, French and Belgian sample data was collected, the analysis  
35  
36 of the tweets content only regarded tweets in English. This might limit the possibility to infer  
37  
38 how the Belgian Twitter users felt as a whole. Third, Belgian Twitter users are not an accurate  
39  
40 representation of the Belgian public, and the tweets must be regarded as an indicative sample.  
41  
42 Fourth, it is not clear, to which extent the sample is subject to automated posting of tweets and  
43  
44 bot activity. These issues are a common concern within the use of Twitter data<sup>32</sup>. Despite its  
45  
46 limitations, the current study adds to the current state of the art Twitter analysis by being the  
47  
48 first to examine the public discourse on Twitter in the first month of the SARS-CoV-2 crises.  
49  
50 As such, it provides important insight into how the Belgian public responded to the far-reaching  
51  
52 measures taken by the government to curtail public life.  
53  
54  
55  
56  
57  
58  
59  
60

**Running head: TWITTER DISCOURSE ON SARS-COV-2****Disclosure Statement**

The authors declare that they have no relevant or material financial interests that relate to the research described in this paper.

**References**

1. McCombs ME, Shaw DL. The Agenda-Setting Function of Mass Media. *Public Opinion Quarterly* 1972; 36:176.
2. McCombs ME, Shaw DL, Weaver DH. New Directions in Agenda-Setting Theory and Research. *Mass Communication and Society* 2014; 17:781–802.
3. Valenzuela S, Puente S, Flores PM. Comparing Disaster News on Twitter and Television: an Intermedia Agenda Setting Perspective. *Journal of Broadcasting & Electronic Media* 2017; 61:615–37.
4. Su Y, Borah P. Who is the agenda setter? Examining the intermedia agenda-setting effect between Twitter and newspapers. *Journal of Information Technology & Politics* 2019; 16:236–49.
5. Vos SC, Buckner MM. Social Media Messages in an Emerging Health Crisis: Tweeting Bird Flu. *Journal of health communication* 2016; 21:301–8.
6. Glowacki EM, Lazard AJ, Wilcox GB, Mackert M, Bernhardt JM. Identifying the public's concerns and the Centers for Disease Control and Prevention's reactions during a health crisis: An analysis of a Zika live Twitter chat. *American journal of infection control* 2016; 44:1709–11.
7. Jain VK, Kumar S. An Effective Approach to Track Levels of Influenza-A (H1N1) Pandemic in India Using Twitter. *Procedia Computer Science* 2015; 70:801–7.

**Running head: TWITTER DISCOURSE ON SARS-COV-2**

- 1  
2  
3 8. Signorini A, Segre AM, Polgreen PM. The use of Twitter to track levels of disease activity  
4 and public concern in the U.S. during the influenza A H1N1 pandemic. PloS one 2011;  
5  
6 6:e19467.  
7  
8
- 9  
10  
11 9. Chan M-PS, Winneg K, Hawkins L, Farhadloo M, Jamieson KH, Albarracín D. Legacy  
12 and social media respectively influence risk perceptions and protective behaviors during  
13 emerging health threats: A multi-wave analysis of communications on Zika virus cases.  
14 Social science & medicine (1982) 2018; 212:50–9.  
15  
16
- 17  
18  
19 10. Chini M. 110 Belgians quarantined in hotel in Tenerife due to coronavirus. Brussels  
20 Times 2020, February 25.  
21  
22
- 23  
24  
25 11. Peltier E, Minder R. Hundreds Confined to Tenerife Hotel for 14 Days Over  
26 Coronavirus Fears. The New York Times 2020, February 26.  
27  
28
- 29  
30  
31 12. Gentry J. (2016) *R Based Twitter Client*.  
32  
33
- 34  
35  
36 13. Feinerer I, Hornik K, Meyer D. Text Mining Infrastructure in R. Journal of Statistical  
37 Software 2008; 25.  
38
- 39  
40  
41 14. Silge J, Robinson D. tidytext: Text Mining and Analysis Using Tidy Data Principles  
42 in R. The Journal of Open Source Software 2016; 1:37.  
43  
44
- 45  
46  
47 15. Porter MF. (2001). Snowball: A language for stemming algorithms. Available at  
48 <http://snowball.tartarus.org/texts/introduction.html>.  
49
- 50  
51  
52 16. Mohammad S, Turney P. (2010) Emotions Evoked by Common Words and Phrases:  
53 Using Mechanical Turk to Create an Emotion Lexicon. In: Inkpen D, Strapparava C, eds.  
54 *Proceedings of the NAACL HLT 2010 Workshop on Computational Approaches to*  
55 *Analysis and Generation of Emotion in Text*. Los Angeles, CA: Association for  
56  
57  
58  
59  
60 Computational Linguistics, pp. 26–34.

**Running head: TWITTER DISCOURSE ON SARS-COV-2**

17. Grün B, Hornik K. topicmodels : An R Package for Fitting Topic Models. *Journal of Statistical Software* 2011; 40.
18. Federal Public Service (FPS) Health, Food Chain Safety and Environment, Belgium. 6 new cases of Covid-19 by the end of the spring holidays. Available at <https://www.info-coronavirus.be/en/news/6-new-cases-of-covid-19-by-the-end-of-the-spring-holidays/>.
19. Gehrke L. Belgium confirms its first coronavirus death. *Politico* 2020, March 11.
20. Chini M. Closing all schools is 'out of the question', says Education Minister. *Brussels Times* 2020, March 12.
21. Chini M, Johnson J. Coronavirus: What are the new measures in Belgium? *Brussels Times* 2020, March 17.
22. Cochrane E, Fandos N. Senate Approves \$2 Trillion Stimulus After Bipartisan Deal. 2020, March 25.
23. Maier D, Waldherr A, Miltner P et al. Applying LDA Topic Modeling in Communication Research: Toward a Valid and Reliable Methodology. *Communication Methods and Measures* 2018; 12:93–118.
24. Arun R, Suresh V, Veni Madhavan CE, Narasimha Murthy MN. (2010) On Finding the Natural Number of Topics with Latent Dirichlet Allocation: Some Observations. In: Zaki MJ, Yu JX, Ravindran B, Pudi V, eds. *Advances in Knowledge Discovery and Data Mining*. Berlin, Heidelberg: Springer Berlin Heidelberg, pp. 391–402.
25. Cao J, Xia T, Li J, Zhang Y, Tang S. A density-based method for adaptive LDA model selection. *Neurocomputing* 2009; 72:1775–81.
26. Griffiths TL, Steyvers M. Finding scientific topics. *Proceedings of the National Academy of Sciences of the United States of America* 2004; 101 Suppl 1:5228–35.



**Running head: TWITTER DISCOURSE ON SARS-COV-2**

- 1  
2  
3 27. Acar A, Muraki Y. Twitter for crisis communication: lessons learned from Japan's  
4 tsunami disaster. *International Journal of Web Based Communities* 2011; 7:392.  
5  
6  
7  
8  
9 28. Terpstra T, Stronkman R, Vries A de, Paradies GL. (2012) Towards a realtime Twitter  
10 Analysis during crises for operational crisis management. In: Rothkrantz L, Ristvej J,  
11 Franco Z, eds. *Proceedings of the 9th International ISCRAM Conference*, pp. 1–9.  
12  
13  
14  
15  
16 29. Wang B, Zhuang J. Crisis information distribution on Twitter: a content analysis of  
17 tweets during Hurricane Sandy. *Natural Hazards* 2017; 89:161–81.  
18  
19  
20  
21  
22 30. Szomszor M, Kostkova P, Louis CS. (2011) Twitter Informatics: Tracking and  
23 Understanding Public Reaction during the 2009 Swine Flu Pandemic. In: *2011*  
24 *IEEE/WIC/ACM International Conferences on Web Intelligence and Intelligent Agent*  
25 *Technology*. IEEE, pp. 320–3.  
26  
27  
28  
29  
30  
31 31. Ahmed W, Bath PA, Sbaffi L, Demartini G. (07182018) Moral Panic through the  
32 Lens of Twitter. In: *Proceedings of the 9th International Conference on Social Media and*  
33 *Society*. New York, NY, USA: ACM, pp. 217–21.  
34  
35  
36  
37  
38  
39 32. Bruns A, Stieglitz S. Twitter data: What do they represent? *it - Information*  
40 *Technology* 2014; 56.  
41  
42  
43  
44  
45  
46

**List of figure legends:**

47  
48  
49  
50 Figure 1: Number of tweets and retweets over time  
51

52  
53 Figure 2: Bigram Network  
54

55  
56 Figure 3: Topic Models  
57

58  
59 Figure 4: Emotions within tweets over time  
60



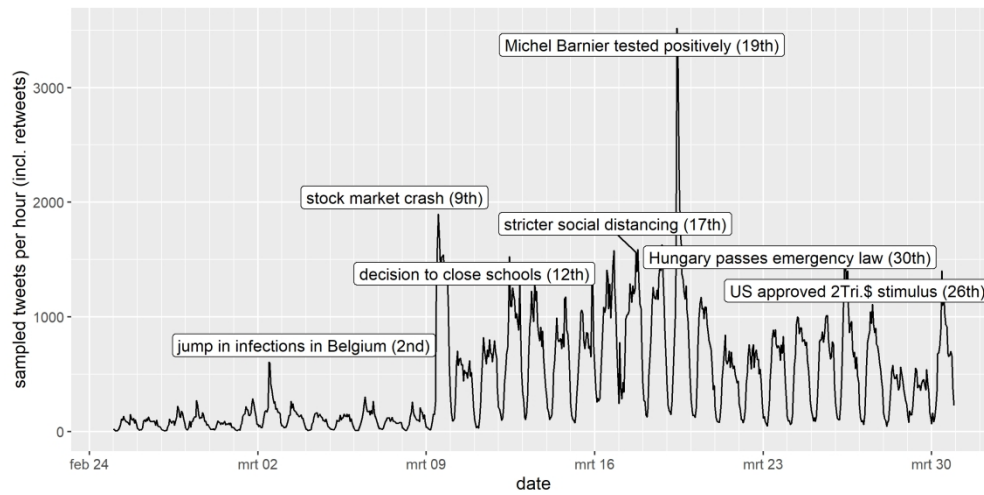


Figure 1: Number of tweets and retweets over time

238x119mm (300 x 300 DPI)



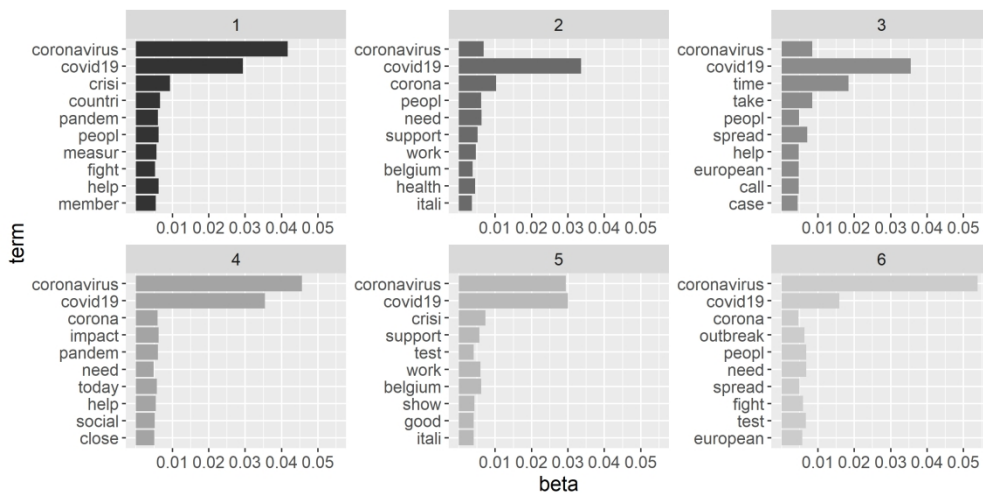


Figure 3: Topic models

238x119mm (300 x 300 DPI)

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

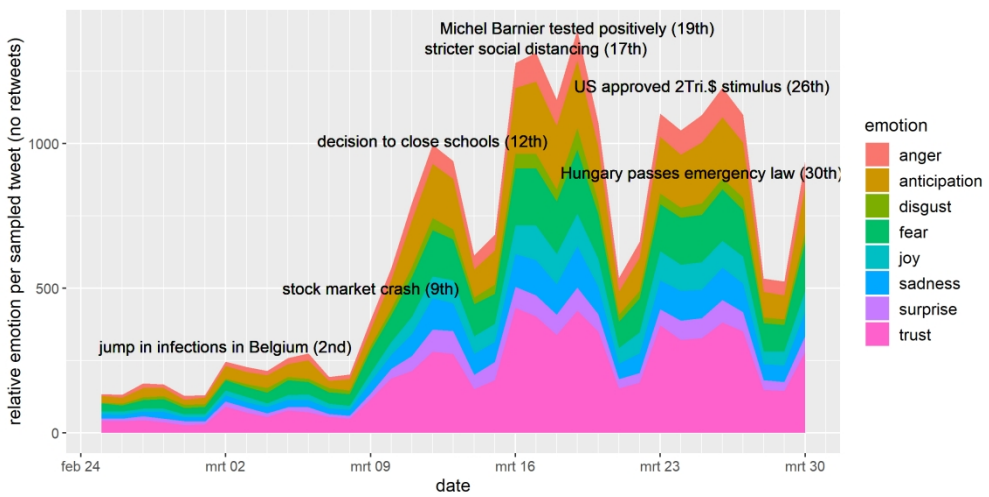


Figure 4: Emotions within tweets over time

238x119mm (300 x 300 DPI)

# #coronavirus: Monitoring the Belgian Twitter Discourse on the Sars-Cov-2 Pandemic

## Technical Appendix

Authors anonymous for peer review

20 July, 2020

### Setup

```
library (twitter)
library (ROAuth)
library (rtweet)
library (openxlsx)
library (dummies)
library (prodlim)
library (ggplot2)
library (dplyr)
library (tm)
library (SnowballC)
library (wordcloud)
library (RColorBrewer)
library (textdata)
library (tidytext)
library (SentimentAnalysis)
library (syuzhet)
library (reshape2)
library (tidyr)
library (igraph)
library (tidygraph)
library (ggraph)
library (topicmodels)
library (scales)
library (doParallel)
library (ldatuning)

#####connect to the API#####

appname <- "corona_extraction"
key <- "XXXXXXXXXXXXXXXXXXXX"
secret <- "XXXXXXXXXXXXXXXXXXXX"
access_token <- "XXXXXXXXXXXXXXXXXXXX"
access_secret <- "XXXXXXXXXXXXXXXXXXXX"
credent <- c(key, secret, access_token, access_secret)

#setup_twitter_oauth(consumer_key = key, consumer_secret = secret,
```

```

1
2
3  access_token = access_token, access_secret = access_secret)
4  #commented out because this script is only used for reporting
5

```

```

6  token <- create_token(
7    app = appname,
8    consumer_key = key,
9    consumer_secret = secret,
10   access_token = access_token,
11   access_secret = access_secret)
12

```

Hashtags to follow:

- coronavid19
- coronavirus
- covid19
- corona
- coronavirusoutbreak
- COVID19BE (from 11th March on)
- coronabelgie (from 11th March on)

We followed these hashtags from the 25th of February to the 31st of March (exclusive the 31st).

## Scraping the Data

```

31 #####set a geocode#####
32

```

```

33 #point coordinates for belgium are

```

```

34 #lat      Lon

```

```

35 #50.640281 4.666715

```

```

36
37 #A radius of 120 KM equals the size of Belgium
38

```

The code below for scraping is only a sample of the scraping code because it was shortened for clarity. The originally used code scraped in smaller timeframes (up to 1 day per language per hashtag per query) to not exceed the limit of the Twitter API which is 18000 tweets per query. The code below shows 15 queries. In total 120 to 150 queries were used.

```

45 twt_coronavid19_NL_02_25 <- twListToDF(searchTwitter("#coronavid19", since =
46 "2020-02-25", until = "2020-03-01", lang = "nl", n = 5000, geocode =
47 "50.640281,4.666715,120km"))
48 twt_coronavid19_NL_02_25$language <- "NL"
49 twt_coronavid19_FR_02_25 <- twListToDF(searchTwitter("#coronavid19", since =
50 "2020-02-25", until = "2020-03-01", lang = "fr", n = 5000, geocode =
51 "50.640281,4.666715,120km"))
52 twt_coronavid19_FR_02_25$language <- "FR"
53 twt_coronavid19_EN_02_25 <- twListToDF(searchTwitter("#coronavid19", since =
54 "2020-02-25", until = "2020-03-01", lang = "en", n = 5000, geocode =
55 "50.640281,4.666715,120km"))
56

```

```
1
2
3 twt_coronavid19_EN_02_25$language <- "EN"
4 twt_coronavid19_02_25 <- rbind (twt_coronavid19_NL_02_25,
5 twt_coronavid19_FR_02_25, twt_coronavid19_EN_02_25)
6 twt_coronavid19_02_25$ht_coronavid19 <- "coronavid19"
7
8
9 twt_coronavirus_NL_02_25 <- twListToDF(searchTwitter("#coronavirus", since =
10 "2020-02-25", until = "2020-03-01", lang = "nl", n = 5000, geocode =
11 "50.640281,4.666715,120km"))
12 twt_coronavirus_NL_02_25$language <- "NL"
13 twt_coronavirus_FR_02_25 <- twListToDF(searchTwitter("#coronavirus", since =
14 "2020-02-25", until = "2020-03-01", lang = "fr", n = 5000, geocode =
15 "50.640281,4.666715,120km"))
16 twt_coronavirus_FR_02_25$language <- "FR"
17 twt_coronavirus_EN_02_25 <- twListToDF(searchTwitter("#coronavirus", since =
18 "2020-02-25", until = "2020-03-01", lang = "en", n = 5000, geocode =
19 "50.640281,4.666715,120km"))
20 twt_coronavirus_EN_02_25$language <- "EN"
21 twt_coronavirus_02_25 <- rbind (twt_coronavirus_NL_02_25,
22 twt_coronavirus_FR_02_25, twt_coronavirus_EN_02_25)
23 twt_coronavirus_02_25$hashtag <- "coronavirus"
24
25
26 twt_covid19_NL_02_25 <- twListToDF(searchTwitter("#covid19", since = "2020-
27 02-25", until = "2020-03-01", lang = "nl", n = 5000, geocode =
28 "50.640281,4.666715,120km"))
29 twt_covid19_NL_02_25$language <- "NL"
30 twt_covid19_FR_02_25 <- twListToDF(searchTwitter("#covid19", since = "2020-
31 02-25", until = "2020-03-01", lang = "fr", n = 5000, geocode =
32 "50.640281,4.666715,120km"))
33 twt_covid19_FR_02_25$language <- "FR"
34 twt_covid19_EN_02_25 <- twListToDF(searchTwitter("#covid19", since = "2020-
35 02-25", until = "2020-03-01", lang = "en", n = 5000, geocode =
36 "50.640281,4.666715,120km"))
37 twt_covid19_EN_02_25$language <- "EN"
38 twt_covid19_02_25 <- rbind (twt_covid19_NL_02_25, twt_covid19_FR_02_25,
39 twt_covid19_EN_02_25)
40 twt_covid19_02_25$hashtag <- "covid19"
41
42
43 twt_corona_NL_02_25 <- twListToDF(searchTwitter("#corona", since = "2020-02-
44 25", until = "2020-03-01", lang = "nl", n = 5000, geocode =
45 "50.640281,4.666715,120km"))
46 twt_corona_NL_02_25$language <- "NL"
47 twt_corona_FR_02_25 <- twListToDF(searchTwitter("#corona", since = "2020-02-
48 25", until = "2020-03-01", lang = "fr", n = 5000, geocode =
49 "50.640281,4.666715,120km"))
50 twt_corona_FR_02_25$language <- "FR"
51 twt_corona_EN_02_25 <- twListToDF(searchTwitter("#corona", since = "2020-02-
52 25", until = "2020-03-01", lang = "en", n = 5000, geocode =
53 "50.640281,4.666715,120km"))
54 twt_corona_EN_02_25$language <- "EN"
55 twt_corona_02_25 <- rbind (twt_corona_NL_02_25, twt_corona_FR_02_25,
```

```

1
2
3 twt_corona_EN_02_25)
4 twt_corona_02_25$hashtag <- "corona"
5
6
7 twt_coronavirusoutbreak_NL_02_25 <-
8 twListToDF(searchTwitter("#coronavirusoutbreak", since = "2020-02-25", until
9 = "2020-03-01", lang = "nl", n = 5000, geocode = "50.640281,4.666715,120km"))
10 twt_coronavirusoutbreak_NL_02_25$language <- "NL"
11 twt_coronavirusoutbreak_FR_02_25 <-
12 twListToDF(searchTwitter("#coronavirusoutbreak", since = "2020-02-25", until
13 = "2020-03-01", lang = "fr", n = 5000, geocode = "50.640281,4.666715,120km"))
14 twt_coronavirusoutbreak_FR_02_25$language <- "FR"
15 twt_coronavirusoutbreak_EN_02_25 <-
16 twListToDF(searchTwitter("#coronavirusoutbreak", since = "2020-02-25", until
17 = "2020-03-01", lang = "en", n = 5000, geocode = "50.640281,4.666715,120km"))
18 twt_coronavirusoutbreak_EN_02_25$language <- "EN"
19 twt_coronavirusoutbreak_02_25 <- rbind(twt_coronavirusoutbreak_NL_02_25,
20 twt_coronavirusoutbreak_FR_02_25, twt_coronavirusoutbreak_EN_02_25)
21 twt_coronavirusoutbreak_02_25$hashtag <- "coronavirusoutbreak"
22
23
24 twt_scrape_02_25 <- rbind(twt_coronavid19_02_25, twt_coronavirus_02_25,
25 twt_covid19_02_25, twt_corona_02_25, twt_coronavirusoutbreak_02_25)
26
27 saveRDS(twt_scrape_02_25, file =
28 "C:\\Users\\u0135880\\Dropbox\\Twitter_Project\\scrape_02_25_to_03_01.rds")
29 write.xlsx(twt_scrape_02_25, file =
30 "C:\\Users\\u0135880\\Dropbox\\Twitter_Project\\scrape_02_25_to_03_01.xlsx")
31
32
33 #####remove duplicated tweets#####
34 #only contains the first 17 columns
35 unique1 <- twt_scrape_02_25[,1:17]
36
37 #dupletes is a variable containing the number of the tweet it is equal to
38 unique1$duplete <- row.match(unique1, unique1)
39 #e.g. if number != duplete it is a duplicated tweet
40 unique1$number <- c(1:nrow(unique1))
41 unique1$hashtag <- twt_scrape_02_25$hashtag
42
43
44 #create hashtag dummy variables
45 unique1 <- cbind(unique1, dummy(unique1$hashtag, sep = "_"))
46
47 #unique_data a subset of unique 1 which only contains the duplicated tweets
48 #sorry for the incoherent name
49 unique_data <- subset(unique1, unique1$duplete != unique1$number)
50
51 #unique2 only contains the original tweets which have duplicates
52 unique2 <- unique1[unique_data$duplete,]
53
54
55 #unique3 contains the first original cases with the added dummies
56 unique3 <- unique2
57
58
59
60

```



```

1
2
3 unique3$unique1_corona <- unique3$unique1_corona+unique_data$unique1_corona
4 unique3$unique1_coronavid19 <-
5 unique3$unique1_coronavid19+unique_data$unique1_coronavid19
6 unique3$unique1_coronavirus <-
7 unique3$unique1_coronavirus+unique_data$unique1_coronavirus
8 unique3$unique1_coronavirusoutbreak <-
9 unique3$unique1_coronavirusoutbreak+unique_data$unique1_coronavirusoutbreak
10 unique3$unique1_covid19 <-
11 unique3$unique1_covid19+unique_data$unique1_covid19
12
13
14 #create new dataframe for further operations
15 unique4 <- unique1
16
17 #replace the values for the dummy variables of the original variables with
18 the dupletes in the full dataset
19 unique4$unique1_corona[match(unique3$number, unique4$number)] <-
20 unique3$unique1_corona
21 unique4$unique1_coronavirus[match(unique3$number, unique4$number)] <-
22 unique3$unique1_coronavirus
23 unique4$unique1_coronavid19[match(unique3$number, unique4$number)] <-
24 unique3$unique1_coronavid19
25 unique4$unique1_coronavirusoutbreak[match(unique3$number, unique4$number)] <-
26 unique3$unique1_coronavirusoutbreak
27 unique4$unique1_covid19[match(unique3$number, unique4$number)] <-
28 unique3$unique1_covid19
29
30
31 #create a subset which does not contain the dupletes
32 unq_twt_02_25 <- subset(unique4, number==duplete)
33
34 #unq_twt_02_25 is now a dataframe with contains the unique tweets and the
35 dummy coded hashtag count
36 unq_twt_02_25$duplete <- NULL
37 unq_twt_02_25$number <- NULL
38 unq_twt_02_25$hashtag <- NULL
39
40
41 #rename the colnums to a more handy name
42 colnames (unq_twt_02_25)[colnames(unq_twt_02_25) == "unique1_corona"] <-
43 "ht_corona"
44 colnames (unq_twt_02_25)[colnames(unq_twt_02_25) == "unique1_coronavid19"] <-
45 "ht_coronavid19"
46 colnames (unq_twt_02_25)[colnames(unq_twt_02_25) == "unique1_coronavirus"] <-
47 "ht_coronavirus"
48 colnames (unq_twt_02_25)[colnames(unq_twt_02_25) ==
49 "unique1_coronavirusoutbreak"] <- "ht_coronavirusoutbreak"
50 colnames (unq_twt_02_25)[colnames(unq_twt_02_25) == "unique1_covid19"] <-
51 "ht_covid19"
52
53
54 saveRDS(unq_twt_02_25, file =
55 "C:\\Users\\u0135880\\Dropbox\\Twitter_Project\\unq_twt_02_25_to_03_01.rds")
56
57
58
59
60

```

```
write.xlsx(unq_twt_02_25, file =
"C:\\Users\\u0135880\\Dropbox\\Twitter_Project\\unq_twt_02_25_to_03_01.xlsx")
```

## Merge Scraped Data and Count Tweets

```
#merging data from the 25th February to the 31st of March
```

```
#Load data and add some NAs because there was a new hashtag regarded since
scrape 03_11
```

```
unq_twt_02_25 <- readRDS (file =
"C:\\Users\\u0135880\\Dropbox\\Twitter_Project\\unq_twt_02_25_to_03_01.rds")
```

```
unq_twt_02_25$ht_COVID19BE <- NA
```

```
unq_twt_02_25$ht_Coronabelgie <- NA
```

```
unq_twt_03_01 <- readRDS (file =
```

```
"C:\\Users\\u0135880\\Dropbox\\Twitter_Project\\unq_twt_03_01_to_03_05.rds")
```

```
unq_twt_03_01$ht_COVID19BE <- NA
```

```
unq_twt_03_01$ht_Coronabelgie <- NA
```

```
unq_twt_03_05 <- readRDS (file =
```

```
"C:\\Users\\u0135880\\Dropbox\\Twitter_Project\\unq_twt_03_05_to_03_11.rds")
```

```
unq_twt_03_05$ht_COVID19BE <- NA
```

```
unq_twt_03_05$ht_Coronabelgie <- NA
```

```
unq_twt_03_11 <- readRDS (file =
```

```
"C:\\Users\\u0135880\\Dropbox\\Twitter_Project\\unq_twt_03_11_to_03_13.rds")
```

```
unq_twt_03_13 <- readRDS (file =
```

```
"C:\\Users\\u0135880\\Dropbox\\Twitter_Project\\unq_twt_03_13_to_03_19.rds")
```

```
unq_twt_03_19 <- readRDS (file =
```

```
"C:\\Users\\u0135880\\Dropbox\\Twitter_Project\\unq_twt_03_19_to_03_23.rds")
```

```
unq_twt_03_23 <- readRDS (file =
```

```
"C:\\Users\\u0135880\\Dropbox\\Twitter_Project\\unq_twt_03_23_to_03_26.rds")
```

```
unq_twt_03_26 <- readRDS (file =
```

```
"C:\\Users\\u0135880\\Dropbox\\Twitter_Project\\unq_twt_03_26_to_03_31.rds")
```

```
twts <- rbind (unq_twt_02_25, unq_twt_03_01, unq_twt_03_05, unq_twt_03_11,
unq_twt_03_13, unq_twt_03_19, unq_twt_03_23, unq_twt_03_26)
```

```
twts_nort <- subset (twts, twts$isRetweet==FALSE)
```

```
twts_nort_EN <- subset (twts_nort, twts_nort$language=="EN")
```

```
twts_nort_FR <- subset (twts_nort, twts_nort$language=="FR")
```

```
twts_nort_NL <- subset (twts_nort, twts_nort$language=="NL")
```

```
#converting
```

```
twts$day_post <- NA
```

```
twts$day_post <- substr (twts$created,1,10)
```

```
twts$day_post <- as.Date (twts$day_post, format = "%Y-%m-%d")
```

```
#count the number of tweets per hour
```

```
twts$hour_post <- NA
```

```
twts$hour_post <- substr (twts$created,1,13)
```

```
twts_per_hour <- twts %>% count(twts$hour_post)
```

```
twts_per_hour$`twts$hour_post` <- as.POSIXct (twts_per_hour$`twts$hour_post`,
format = "%Y-%m-%d %H") #tweets within the hour
```

```

1
2
3
4 #count the numbers of tweets per day
5 twts_per_day <- twts %>% count(twts$day_post)
6 twts_per_day$`twts$day_post` <- as.Date (twts_per_day$`twts$day_post`, format
7 = "%Y-%m-%d")
8
9

```

## Data Cleaning and Stemming

```

10 docs <- Corpus(VectorSource(twts_nort_EN$text))
11
12
13 #####cleaning#####
14
15 toSpace <- content_transformer(function (x , pattern ) gsub(pattern, " ", x))
16 docs <- tm_map(docs, toSpace, "/" )
17 docs <- tm_map(docs, toSpace, "@")
18 docs <- tm_map(docs, toSpace, "\\|")
19 docs <- tm_map(docs, toSpace, "...")
20
21
22 docs[[27]]$content
23
24 ## [1] "Nasa images show China pollution clear amid slowdown due to
25 #coronavirus. https: t.co esQKoahPjH"
26
27 docs <- tm_map(docs, stripWhitespace)
28 docs <- tm_map(docs, removeWords, c("https", "t.co", "tco", "..."))
29 docs <- tm_map(docs, stripWhitespace)
30 docs <- tm_map(docs, removeWords, c("https", "t.co", "tco", "..."))
31 docs <- tm_map(docs, stripWhitespace)
32 docs <- tm_map(docs, removeWords, stopwords("english"))
33 docs <- tm_map(docs, stripWhitespace)
34 docs <- tm_map(docs, removePunctuation)
35 docs <- tm_map(docs, stripWhitespace)
36 docs <- tm_map(docs, content_transformer(tolower))
37 docs <- tm_map(docs, stemDocument)
38 docs <- tm_map(docs, stripWhitespace)
39
40
41 docs[[27]]$content
42
43 ## [1] "nasa imag show china pollut clear amid slowdown coronavirus
44 esqkoahpjh"
45
46 dtm <- TermDocumentMatrix(docs)
47

```

## Emotions over Time

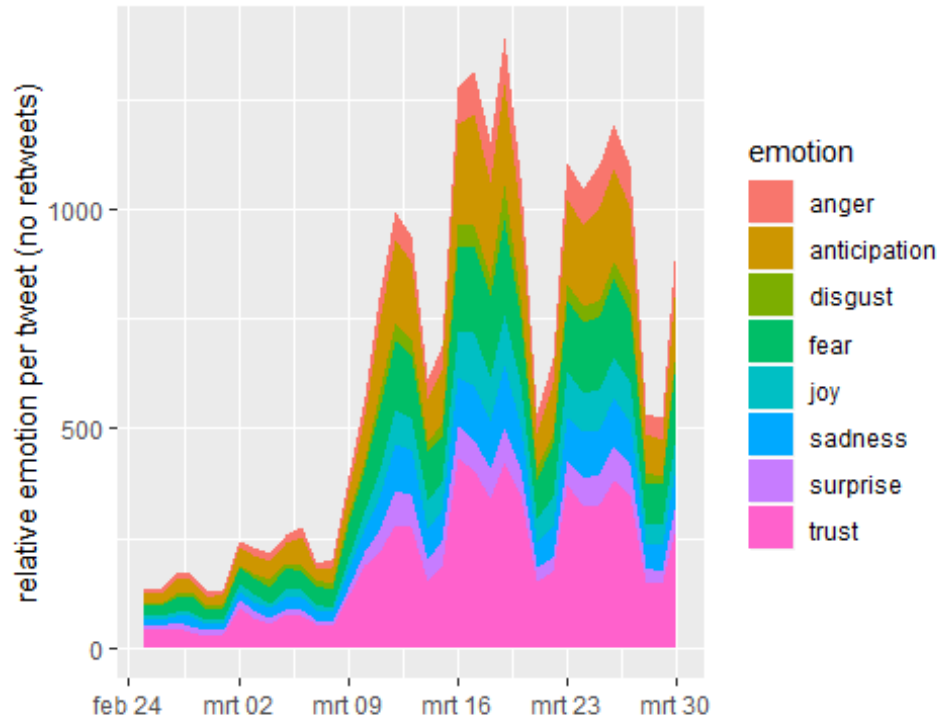
```

48
49 twts_nort_EN_sent <- twts_nort_EN
50
51 #get the emotions
52 sent2 <- get_nrc_sentiment(twts_nort_EN$text)
53
54
55 #bind it back together
56
57
58
59
60

```

```
1
2
3 twts_nort_EN_sent <- cbind (twts_nort_EN_sent, sent2)
4
5 saveRDS(twts_nort_EN_sent, file =
6 "C:\\Users\\u0135880\\Dropbox\\Twitter_Project\\twts_nort_EN_sent.rds")
7
8 twts_nort_EN_sent <- readRDS(file =
9 "C:\\Users\\u0135880\\Dropbox\\Twitter_Project\\twts_nort_EN_sent.rds")
10
11 #####"absolute" timeseries####
12
13 attach(twts_nort_EN_sent)
14
15
16 #create relative count of emotions and sentiments per tweet
17 twts_nort_EN_sent$emotions_sum <- NA
18 twts_nort_EN_sent$emotions_sum <- twts_nort_EN_sent$anger +
19 twts_nort_EN_sent$anticipation + twts_nort_EN_sent$disgust +
20 twts_nort_EN_sent$fear + twts_nort_EN_sent$joy + twts_nort_EN_sent$sadness +
21 twts_nort_EN_sent$surprise + twts_nort_EN_sent$trust
22 #twts_nort_EN_sent$emotions_sum <- anger + anticipation + disgust + fear +
23 joy + sadness + surprise + trust
24
25
26 twts_nort_EN_sent$sent_sum <- negative + positive
27
28 attach(twts_nort_EN_sent)
29
30 twts_nort_EN_sent$anger_rel <- anger / emotions_sum
31 twts_nort_EN_sent$anticipation_rel <- anticipation / emotions_sum
32 twts_nort_EN_sent$disgust_rel <- disgust / emotions_sum
33 twts_nort_EN_sent$fear_rel <- fear / emotions_sum
34 twts_nort_EN_sent$joy_rel <- joy / emotions_sum
35 twts_nort_EN_sent$sadness_rel <- sadness / emotions_sum
36 twts_nort_EN_sent$surprise_rel <- surprise / emotions_sum
37 twts_nort_EN_sent$trust_rel <- trust / emotions_sum
38
39
40 twts_nort_EN_sent$positive_rel <- positive / sent_sum
41 twts_nort_EN_sent$negative_rel <- negative / sent_sum
42
43 #replace nan with 0 (they occurred because the code above divided by 0
44 sometimes.)
45 is.nan.data.frame <- function(x)
46 do.call(cbind, lapply(x, is.nan))
47 twts_nort_EN_sent[is.nan(twts_nort_EN_sent)] <- 0
48
49
50 twts_nort_EN_sent$one <- 1
51
52 #maybe not necessary to convert it, because it is already in kind of a date
53 format.
54 twts_nort_EN_sent$day_post <- NA
55 twts_nort_EN_sent$day_post <- substr (twts_nort_EN_sent$created,1,10)
```

```
1
2
3 twts_nort_EN_sent$day_post <- as.Date (twts_nort_EN_sent$day_post, format =
4 "%Y-%m-%d")
5
6
7 attach(twts_nort_EN_sent)
8 #calculate the sum of emotions per day
9 anger_pd <- as.data.frame(tapply (anger_rel, day_post, sum))
10 anticipation_pd <- as.data.frame(tapply (anticipation_rel, day_post, sum))
11 disgust_pd <- as.data.frame(tapply (disgust_rel, day_post, sum))
12 fear_pd <- as.data.frame(tapply (fear_rel, day_post, sum))
13 joy_pd <- as.data.frame(tapply (joy_rel, day_post, sum))
14 sadness_pd <- as.data.frame(tapply (sadness_rel, day_post, sum))
15 surprise_pd <- as.data.frame(tapply (surprise_rel, day_post, sum))
16 trust_pd <- as.data.frame(tapply (trust_rel, day_post, sum))
17 #calculate the amount of negativity and positivity per day
18 negative_pd <- as.data.frame(tapply (negative_rel, day_post, sum))
19 positive_pd <- as.data.frame(tapply (positive_rel, day_post, sum))
20
21 emotions_pd <- cbind (anger_pd, anticipation_pd, disgust_pd, fear_pd, joy_pd,
22 sadness_pd, surprise_pd, trust_pd)
23 emotions_pd$day_post <- rownames(emotions_pd)
24 emotions_pd$day_post <- as.Date (emotions_pd$day_post, format = "%Y-%m-%d")
25 colnames(emotions_pd) <- c("anger", "anticipation", "disgust", "fear", "joy",
26 "sadness", "surprise", "trust", "day_post")
27
28
29 emotions_plot <- pivot_longer(emotions_pd, ~"day_post", names_to = "emotion")
30
31
32 ggplot(data = emotions_plot, aes(x=day_post, fill=emotion, y = value)) +
33   geom_area() +
34   scale_x_date(breaks = "week", date_labels = "%b %d") +
35   xlab ("") +
36   ylab ("relative emotion per tweet (no retweets)")
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
```

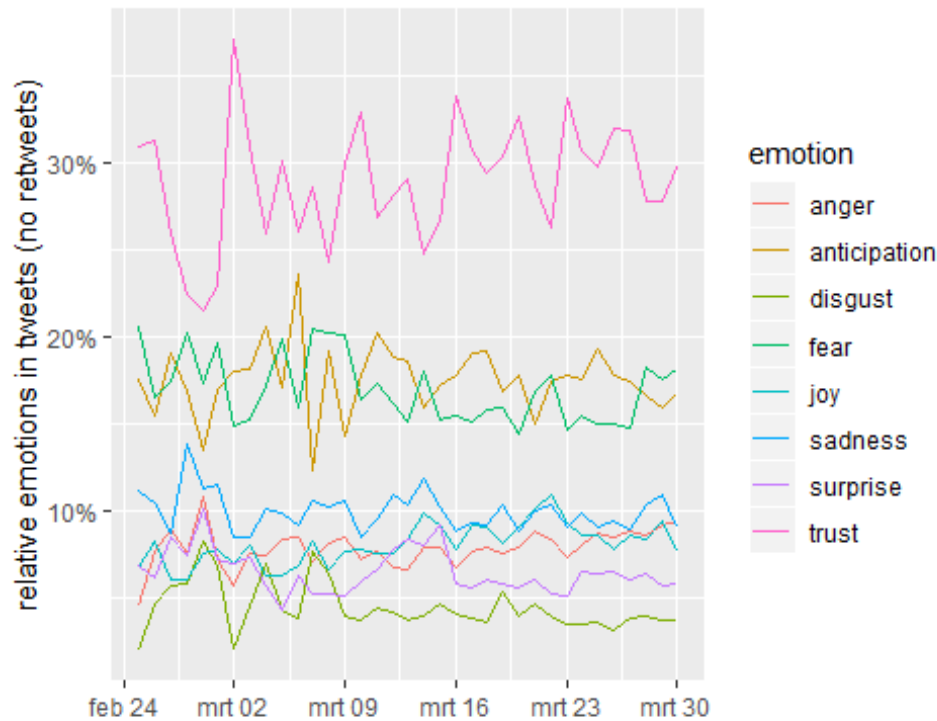


```
#####relative timeseries#####
```

```
emotions_pd_calc <- emotions_pd
#create rowsums of emotions
emotions_pd_calc$rowsum <- rowSums(emotions_pd_calc[,1:8])
#divide values by rowsums to obtain relative amount of emotion on unique
tweets within a given day
emotions_pd_rel <- sweep(emotions_pd_calc[,1:8], 1, emotions_pd_calc$rowsum,
FUN = '/')
emotions_pd_rel$day_post <- emotions_pd$day_post

emotions_rel_plot <- pivot_longer(emotions_pd_rel, ~"day_post", names_to =
"emotion")

ggplot(data = emotions_rel_plot, aes(x=day_post, fill=emotion, y = value,
color = emotion)) +
  geom_line() +
  scale_x_date(breaks = "week", date_labels = "%b %d") +
  scale_y_continuous(labels = scales::percent) +
  ylab("relative emotions in tweets (no retweets)") +
  xlab ("")
```



## Tidy Word Analysis

*#inspired by: <https://www.tidytextmining.com/ngrams.html>*

*#get the cleaned text from the corpus object used before*

```
l <- length(docs)
```

```
twts_EN_nort_clean_list <- lapply(docs[1:l], as.character)
```

```
twts_EN_nort_clean <- unlist(twts_EN_nort_clean_list)
```

```
#####
```

```
twts_tidy_EN <- data.frame(txt = twts_EN_nort_clean, stringsAsFactors = FALSE)
```

*#frequency count of words*

```
twts_tidy_EN %>%
```

```
  unnest_tokens(output = word, input = txt) %>%
```

```
  count(word, sort = TRUE)
```

```
## # A tibble: 69,271 x 2
```

```
##   word      n
```

```
##   <chr>    <int>
```

```
## 1 coronavirus 9274
```

```
## 2 covid19    8991
```

```
## 3 will       1887
```

```
## 4 peopl     1530
```

```
## 5 time      1441
```

```
## 6 corona    1284
```



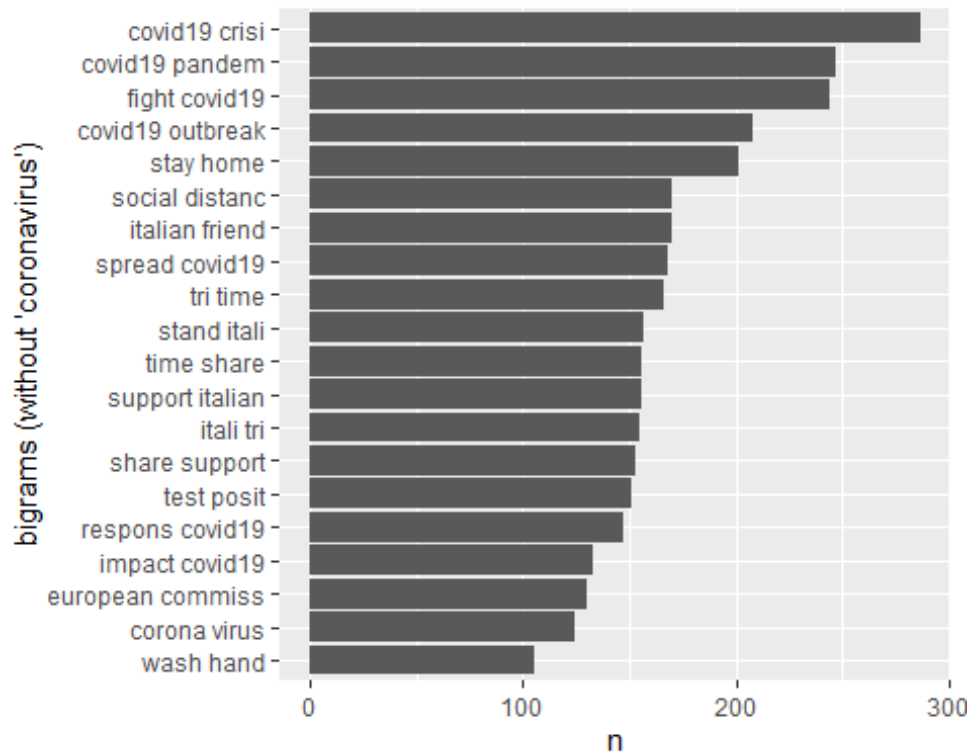




```

1
2
3 bigrams$bigram <- gsub("[^[:alnum:]]", " ", bigrams$bigram) #get rid of all
4 the smilies and other non-alphanumeric characters
5 big_clean <- bigrams[-grep(" ", bigrams$bigram),] #get rid of bigrams that
6 contained single non-alphanumeric characters (e.g. "-covid19")
7 big_clean <- as_tibble(big_clean)
8 head(big_clean, 10)
9
10
11 ## # A tibble: 10 x 2
12 ##   bigram          n
13 ##   <chr>          <int>
14 ## 1 coronavirus covid19    457
15 ## 2 covid19 coronavirus    356
16 ## 3 covid19 crisi          287
17 ## 4 covid19 pandem        247
18 ## 5 fight covid19          244
19 ## 6 member state           241
20 ## 7 coronavirus crisi       237
21 ## 8 coronavirus outbreak    211
22 ## 9 covid19 outbreak        208
23 ## 10 stay home              201
24
25 bigrams_separated <- big_clean %>%
26   separate(bigram, c("word1", "word2"), sep = " ")
27
28 #filter for stop_words
29 bigrams_filtered <- bigrams_separated %>%
30   filter(!word1 %in% stop_words$word) %>%
31   filter(!word2 %in% stop_words$word)
32
33
34 #create a vector which contains both words again
35 bigrams_filtered$words <- paste (bigrams_filtered$word1,
36 bigrams_filtered$word2, sep = " ", collapse = NULL)
37 bigrams_filtered_stripped <- bigrams_filtered[-grep("coronavirus",
38 bigrams_filtered$words),] #get rid of bigrams that contain the term
39
40
41 ggplot(bigrams_filtered_stripped[1:20,], aes(x=reorder(words, n), y=n)) +
42   geom_bar(stat="identity") +
43   coord_flip() +
44   xlab("bigrams (without 'coronavirus')")
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

```



```
#####visualazing bigram plots#####
```

```
# filter for only relatively common combinations
```

```
bigram_graph <- bigrams_filtered %>%
```

```
  filter(n > 20) %>%
```

```
  graph_from_data_frame()
```

```
set.seed(2020)
```

```
ggraph(bigram_graph, layout = "fr") +
```

```
  geom_edge_link(aes(edge_alpha = n), show.legend = FALSE,
```

```
                 end_cap = circle(.07, 'inches')) +
```

```
  geom_node_point(color = "grey", size = 5) +
```

```
  geom_node_text(aes(label = name), vjust = 1, hjust = 1) +
```

```
  theme_void()
```



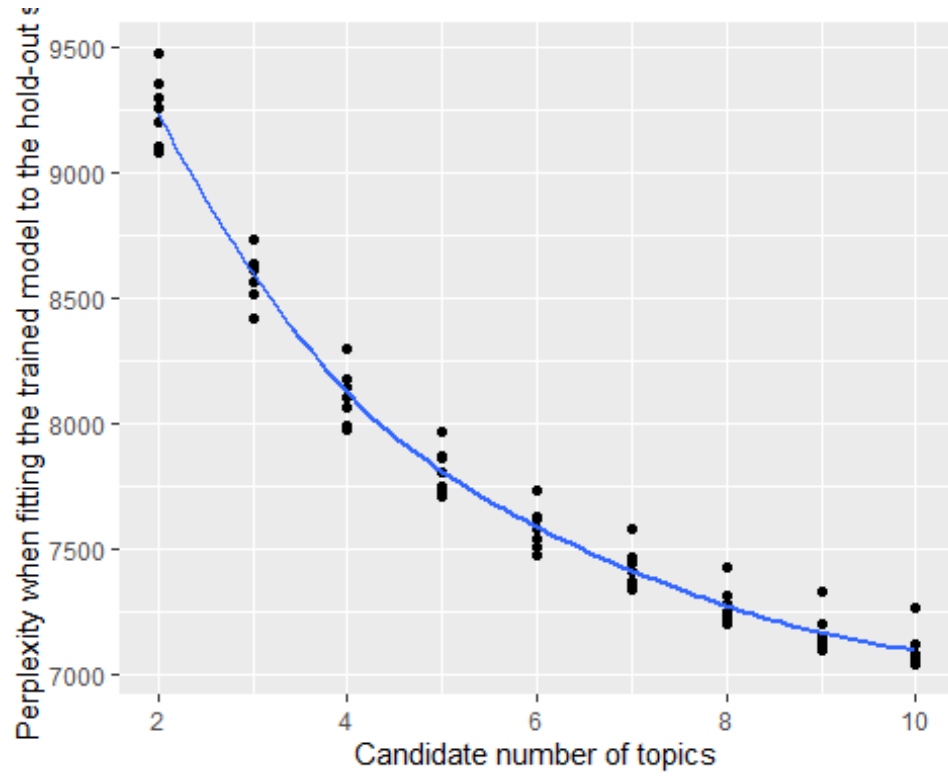
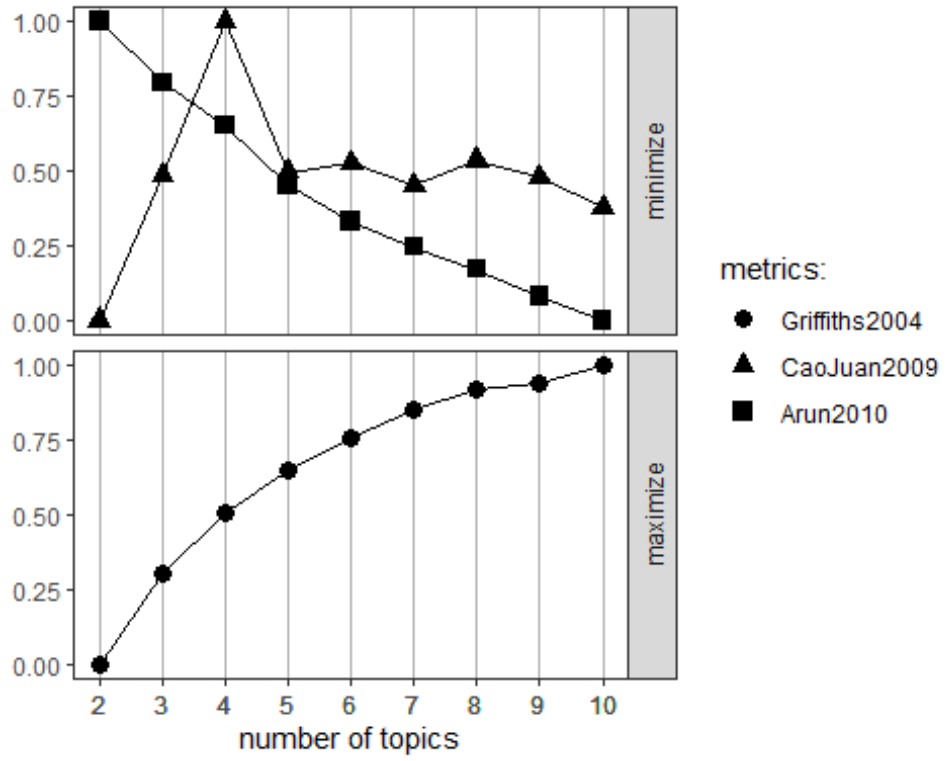
```

1
2
3   arrange(topic, -beta)
4   ap_plot <- ap_top_terms %>%
5     mutate(term = reorder(term, beta))
6
7
8   saverDS(ap_plot, file =
9     "C:\\Users\\u0135880\\Dropbox\\Twitter_Project\\ap_plot.rds")
10
11   #####cross validation of topic models#####
12   #inspired by: https://rpubs.com/MNidhi/NumberoftopicsLDA
13
14   system.time({
15     tunes <- FindTopicsNumber(
16       dtm = dtm2,
17       topics = c(2:10),
18       metrics = c("Griffiths2004", "CaoJuan2009", "Arun2010"),
19       method = "Gibbs",
20       control = list(seed = 77),
21       mc.cores = 4L,
22       verbose = TRUE
23     )
24   })
25
26
27   saverDS(tunes, file =
28     "C:\\Users\\u0135880\\Dropbox\\Twitter_Project\\tunes.rds")
29
30   #key fold cross validation
31
32
33   topics <- c(2:15)
34   burnin = 100
35   iter = 1000
36   keep = 50
37   folds <- 10
38   splitfolds <- sample(1:folds, 23, replace = TRUE)
39   candidate_k <- c(2:10) # candidates for how many topics
40
41   system.time({
42     results <- foreach(j = 1:length(candidate_k), .combine = rbind) %dopar%{
43       k <- candidate_k[j]
44       results_1k <- matrix(0, nrow = folds, ncol = 2)
45       colnames(results_1k) <- c("k", "perplexity")
46       for(i in 1:folds){
47         train_set <- dtm2[splitfolds != i , ]
48         valid_set <- dtm2[splitfolds == i, ]
49
50
51         fitted <- LDA(train_set, k = k, method = "Gibbs",
52                       control = list(burnin = burnin, iter = iter, keep = keep)
53       )
54       results_1k[i,] <- c(k, perplexity(fitted, newdata = valid_set))
55     }
56
57
58
59
60

```

```
1  
2  
3     return(results_1k)  
4   }  
5 })  
6  
7 results_df_10fold_cv <- as.data.frame(results)  
8 saveRDS(results_df_10fold_cv, file =  
9 "C:\\Users\\u0135880\\Dropbox\\Twitter_Project\\results_df_10fold_cv.rds")
```

11 Plots to access the appropriate number of topics:  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

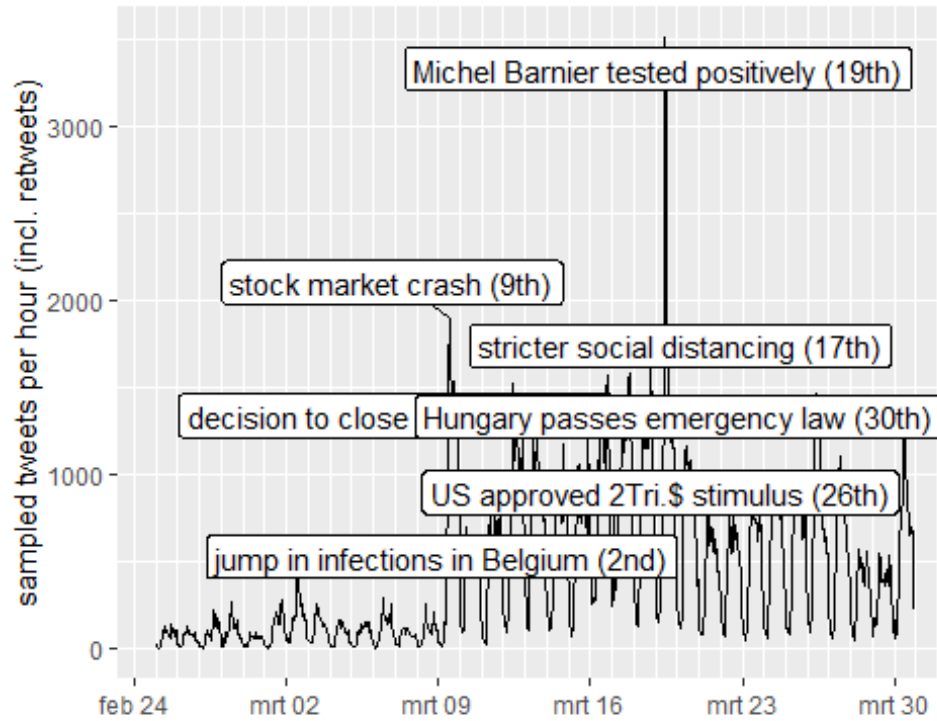


## Reporting for the Paper

```
#####tweets per hour plot

#add some labels
twts_per_hour$label <- NA
twts_per_hour$label[156] <- c("jump in infections in Belgium (2nd)") #March
2nd
twts_per_hour$label[325] <- c("stock market crash (9th)") #March 9th
twts_per_hour$label[396] <- c("decision to close schools (12th)") #March 12th
twts_per_hour$label[522] <- c("stricter social distancing (17th)") #March
17th
twts_per_hour$label[563] <- c("Michel Barnier tested positively (19th)")
#March 19th
twts_per_hour$label[729] <- c("US approved 2Tri.$ stimulus (26th)") #March
26th
twts_per_hour$label[828] <- c("Hungary passes emergency law (30th)") #March
30th

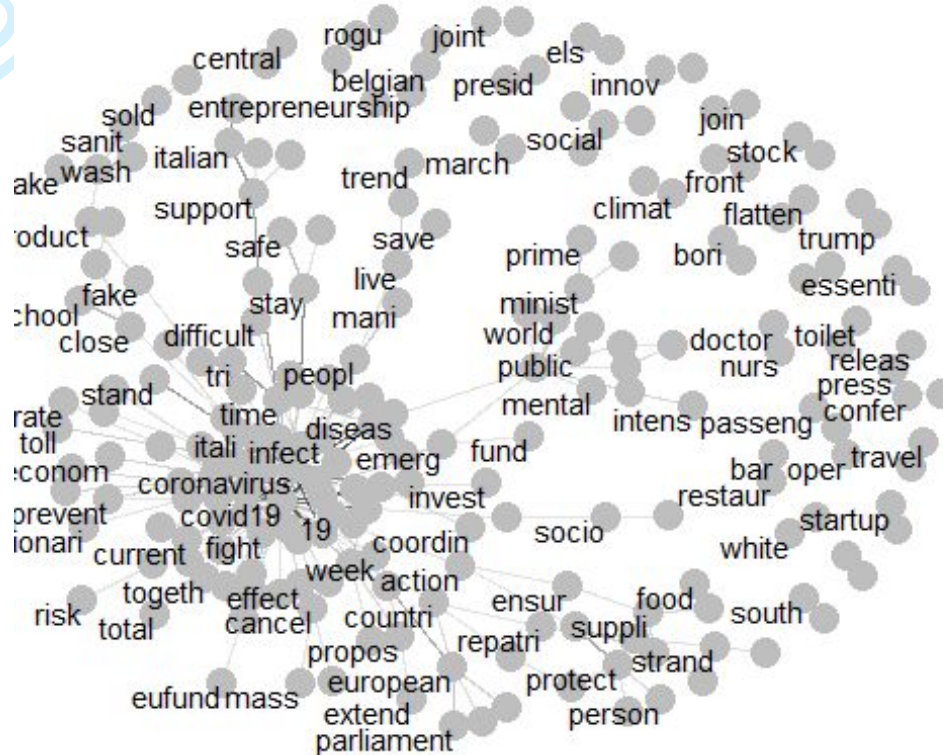
library (ggrepel)
ggplot(twts_per_hour, aes(x=`twts$hour_post`, y=n)) +
  geom_line() +
  scale_x_datetime(date_minor_breaks = "day", date_labels = "%b %d",
date_breaks = "week") +
  ylab ("sampled tweets per hour (incl. retweets)") +
  xlab ("") +
  geom_label_repel (aes (label = label))
```



```
#####bigram network plot#####
```

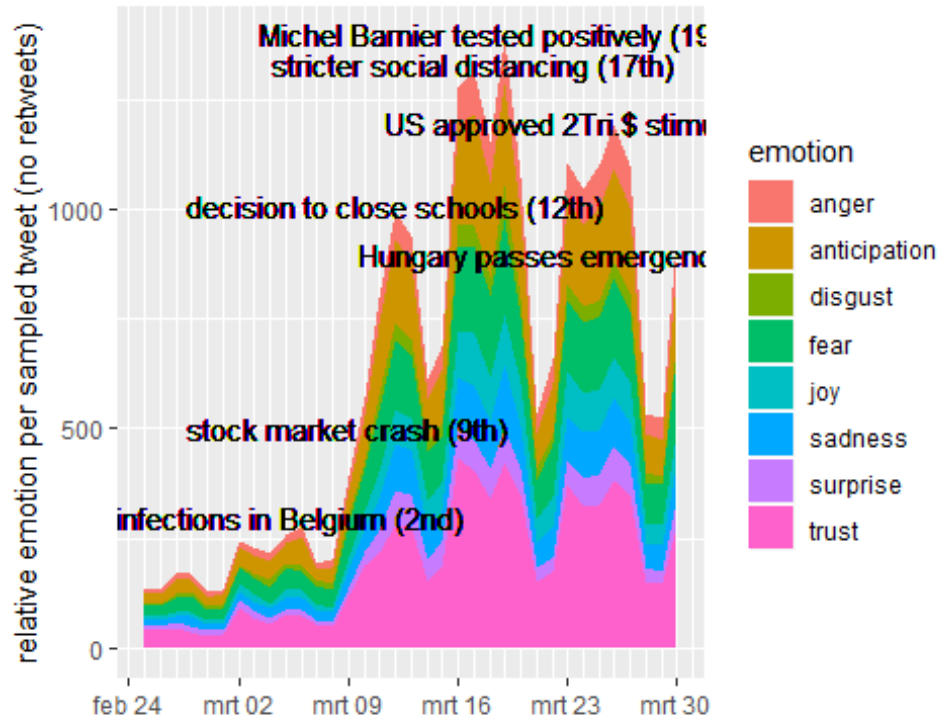
```
set.seed(2020)
ggraph(bigram_graph, layout = "fr") +
  geom_edge_link(aes(edge_alpha = n), show.legend = FALSE,
                end_cap = circle(.07, 'inches')) +
  geom_node_point(color = "grey", size = 5) +
  geom_node_text(aes(label = name), vjust = 1, hjust = 1, check_overlap =
TRUE) + #overlapping words are not plotted to increase clarity
  theme_void()
```





```
#####emotions plot#####
```

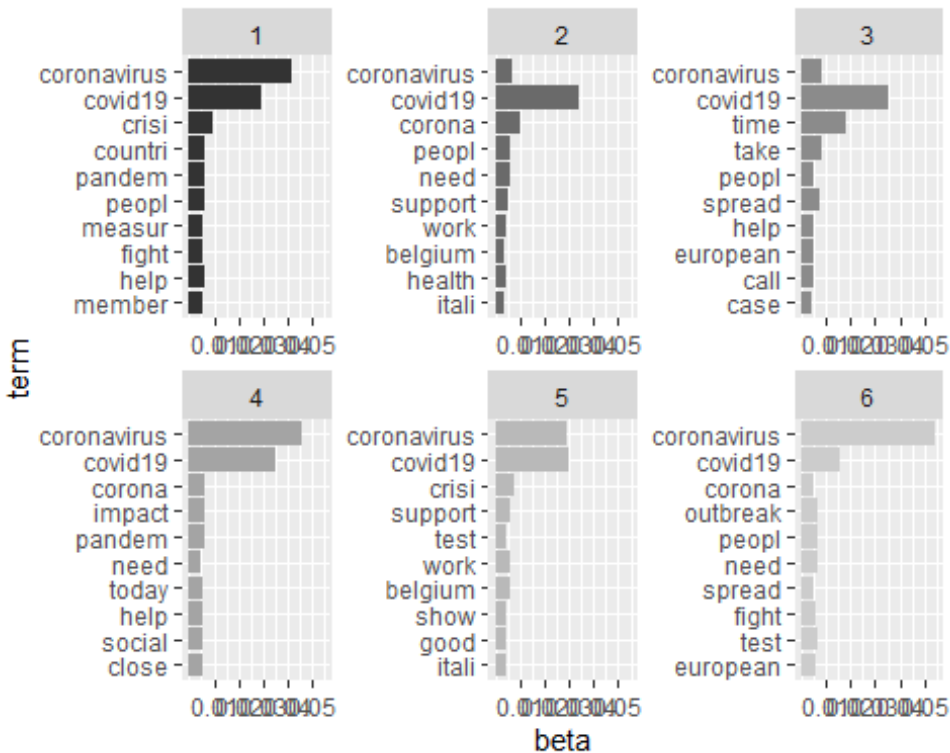
```
ggplot(data = emotions_plot, aes(x=day_post, fill=emotion, y = value)) +
  geom_area() +
  scale_x_date(date_minor_breaks = "day", date_labels = "%b %d", date_breaks
= "week") +
  xlab ("") +
  ylab ("relative emotion per sampled tweet (no retweets)") +
  geom_text(x=as.Date("2020-03-02", format = "%Y-%m-%d"), y=300,
label=c("jump in infections in Belgium (2nd)", color="black") +
  geom_text(x=as.Date("2020-03-09", format = "%Y-%m-%d"), y=500,
label=c("stock market crash (9th)", color="black") +
  geom_text(x=as.Date("2020-03-12", format = "%Y-%m-%d"), y=1010,
label=c("decision to close schools (12th)", color="black") +
  geom_text(x=as.Date("2020-03-17", format = "%Y-%m-%d"), y=1330,
label=c("stricter social distancing (17th)", color="black") +
  geom_text(x=as.Date("2020-03-19", format = "%Y-%m-%d"), y=1400,
label=c("Michel Barnier tested positively (19th)", color="black") +
  geom_text(x=as.Date("2020-03-26", format = "%Y-%m-%d"), y=1200,
label=c("US approved 2Tri.$ stimulus (26th)", color="black") +
  geom_text(x=as.Date("2020-03-26", format = "%Y-%m-%d"), y=900,
label=c("Hungary passes emergency law (30th)", color="black") #X value was
changed to make the text fit on the plot
```



```
#####topic models#####
```

```
ap_plot <- readRDS(file =
"C:\\Users\\u0135880\\Dropbox\\Twitter_Project\\ap_plot.rds")

ggplot(data = ap_plot, aes(term, beta, fill = factor(topic))) +
geom_col(show.legend = FALSE) +
facet_wrap(~ topic, scales = "free") +
coord_flip() +
scale_y_continuous(breaks=c(.01,.02,.03,.04,.05), limits = c(0,0.055)) +
scale_fill_grey()
```



#####reporting for paper#####

*#number of tweets and retweets:*

`nrow(twts)`

## [1] 373908

*#number of original tweets:*

`nrow(twts_nort)`

## [1] 71481

*#number of original tweets in English:*

`nrow(twts_nort_EN)`

## [1] 31454

*#number of original tweets in Dutch:*

`nrow(twts_nort_NL)`

## [1] 24091

*#number of original tweets in French:*

`nrow(twts_nort_FR)`

## [1] 15936

*#relative overall amount of emotions*

`tapply (emotions_rel_plot$value, emotions_rel_plot$emotion, mean)`

```
1
2
3 ##      anger anticipation      disgust      fear      joy
4 sadness
5 ## 0.07824566 0.17502977 0.04393404 0.16957330 0.08018855
6 0.09948661
7 ##      surprise      trust
8 ## 0.06430995 0.28923212
9
10 #all emotions
11 emotions_rel_plot
12
13 ## # A tibble: 280 x 3
14 ##   day_post  emotion      value
15 ##   <date>   <chr>      <dbl>
16 ## 1 2020-02-25 anger      0.0452
17 ## 2 2020-02-25 anticipation 0.175
18 ## 3 2020-02-25 disgust    0.0193
19 ## 4 2020-02-25 fear      0.205
20 ## 5 2020-02-25 joy      0.0671
21 ## 6 2020-02-25 sadness   0.111
22 ## 7 2020-02-25 surprise  0.0677
23 ## 8 2020-02-25 trust     0.309
24 ## 9 2020-02-26 anger     0.0762
25 ## 10 2020-02-26 anticipation 0.154
26 ## # ... with 270 more rows
```