

UNIVERSITA DEGLI STUDI DI PADOVA

DEPARTMENT OF STATISTICAL SCIENCE

Analysis of the topics discussed on Twitter during Covid19 pandemic by Italian politicians

Project for Statistics Method for Big Data

Authors

Boscarini Alessandro, Cifelli Lorenzo, Zampieri Daniele

2020/2021

Preprocessing

Descriptive analysis

[illegible]

Then we analyzed the wordclouds for some politicians. In Figure 2 there is the wordcloud of Matteo Salvini: the most frequent words are about social communication and covid19.

The first is to group the tweets of the same politician in one document and then fit the LDA on this corpus. In this way we obtained longer documents and eliminated the correlation between the tweets of the same politician. However, this approach has not been effective because the corpus of document is small and they are similar to each other.

The second approach is to apply some extension of LDA. One is called Author Topic Model and take account also of the author of the documents. This model provides a three.step document generations: the author is selected, then given the author the topics are chosen and finally, given a topic, the words of the document are generated. However, this model did not bring to satisfactory results, probably because the model is designed for a corpus of documents where every document has many authors. The second extension of LDA considered is Dirichlet Multinomial Mixture model. It assume that every document has at least one topic. This model seems to fit our data very well since the tweets are short text documents and generally are about one topic. Unfortunately we didn't find an implementation in Python.

In conclusion, we decided to present the results of the initial approach that consider every tweet as single independent document because, in our opinion, it provides the best results. However, we considered only the tweets with more than 200 words in order to have fairly long documents.

For the implementation of LDA we us the library Gensim in Python. The estimate of the model is based on variational bayes.

Result

The model needs to set the total number of topics. We fit different model for values ranging from 5 to 50. We chose the best model by looking at which one assigned the topic in most appropriate way. We chose the model with 20 topic. I remind you that the output of LDA is a probability distribution of topics for every tweet and a probability distribution of words for every topic. For some topics (1), the probability distribution of words define a specific themes. For others (2) it is more difficult to define a title for the topic from the probability distribution of words.

Then we defined a title to each topic and we assigned to each tweet the topic with the highest estimated probability thereby to highlight the most discussed topics. From 6 it is possible to see that the most common topics reflect the different aspect of the pandemic. The theme of employment is one of the must discussed which is related to the political debate on that has

Topic 4 aid to enterprises	Topic 6 covid19 emergency	Topic 7 economy	Topic 14 hospital	Topic 16 Europe and Mes
imprese	coronavirus	euro	medici	mes
liquidità	emergenza	miliardi	infermieri	ue
famiglie	governo	stato	lavoro	europa
aziende	misure	risorse	personale	bce
laboratori	covid19	pil	operatori	eurogruppo

Table 1: Consistent topics

Topic 1 intensive care unit	Topic 10 economic crisis	Topic 11 infections and swaps	Topic 18 economy school	Topic 8 undefinable
storia	sicurezza	contagi	crescita	<i>staseraitalia</i>
terapia	p.iva	tv	crescita	vita
intensiva	Draghi	social	investimenti	Bergamo
parlamento	crisi	tamponi	scuole	persone
leggere	misure	video	contagi	sotto

Table 2: Inconsistent topics

taken place on aid to companies, layoffs and VAT. The other themes highlighted are about the government and its decisions, the health emergency, the economic crisis and European Union.

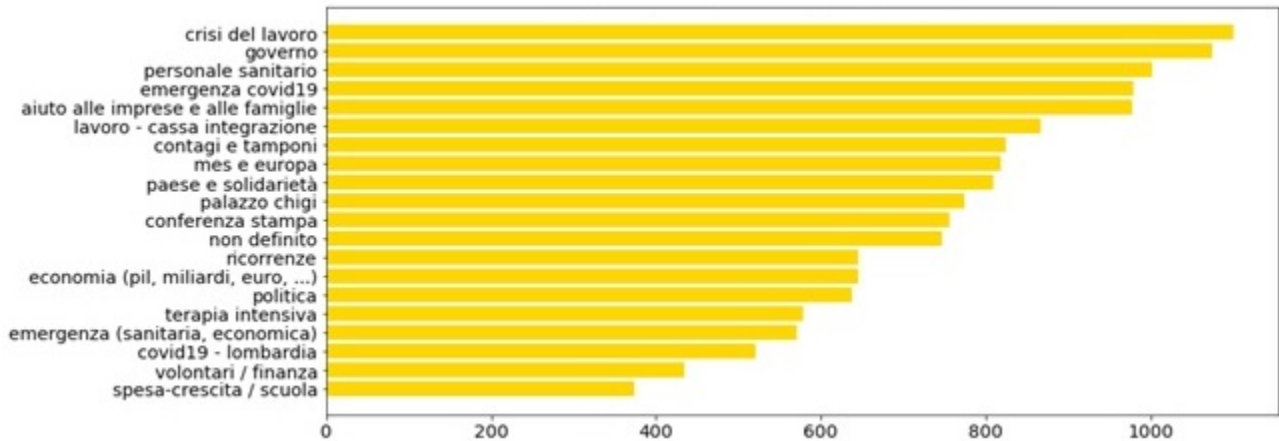


Figure 6: Overall frequency of topics

Finally, we analyzed the frequency of the topics for some politicians (7). For example, Giuseppe Conte (first minister, Movimento 5 Stelle) discussed more about the covid19 emergency. Nicola Zingaretti (Partito Democratico) dedicated more tweets to healthcare workers while ALberto Bagnai (Lega) presents more tweets about European Union and Mes.

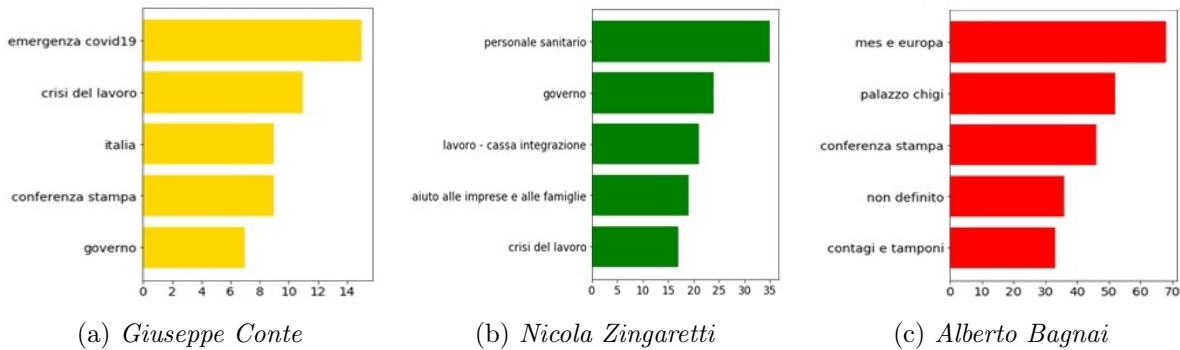


Figure 7: Frequency of topics for politician