

Capstone Project One

Topic modeling for news articles

Milestone Report

by Yiding Weng
April 20, 2019

Problem

With the rapid growth of internet and online information services, printing medias, more and more information is available and accessible. Explosion of information has caused a well recognized information overload problem. There is no time to read everything and yet we have to make critical decisions based on whatever information is available. Without absorbing the text information line by line, how to know what is the significant writing worth paying attention? What is the central theme?

One of the primary applications of natural language processing is to automatically extract what topics people are discussing from large volumes of text. Knowing what people are talking about and understanding their problems and opinions is highly valuable to businesses, administrators, political campaigns. It will be very helpful to have an automated algorithm that can read through the text documents and automatically output the topics discussed.

I will use news articles as the source of information, and extract the naturally discussed topics.

Data

Proposed Data Collection:

All the news: 143,000 articles from 15 American publications

Data Location: <https://www.kaggle.com/snapcrack/all-the-news>

Approach

I will be using the Latent Dirichlet Allocation (LDA) from Gensim package along with the Mallet's implementation (via Gensim). Mallet has an efficient implementation of the LDA. It is known to run faster and gives better topics segregation.

We will also extract the volume and percentage contribution of each topic to get an idea of how important a topic is.

Data wrangling

Download nltk stop words and spacy model

We will need the stop words from NLTK and spacy's en model for text pre-processing. Then we will be using the spacy model for lemmatization. Lemmatization is converting a word to its root word.

For example: the lemma of the word 'machines' is 'machine'. Likewise, 'walking' → 'walk', 'mice' → 'mouse' and so on.

Data sampling

This version of the dataset contains about 142k news articles from 15 different medias. This is available as <https://www.kaggle.com/jannesklaas/analyzing-the-news/data>. Due to the limited computing power, I will take a sample of 30% of the total data for our study in this project.

Data cleaning

'title' and 'content' are the important features that have the text information we are interested to study. Since this title dataset does not have empty fields in 'title' or 'content' information, we will not remove any rows of data.

```
id          0
title       0
publication 0
author      4872
date        798
year        798
month       798
url         16934
content     0
dtype: int64
```

Remove emails and newline characters

There are many emails sign (@), newline and extra spaces that is quite distracting. So get rid of them using regular expressions. After removing the emails and extra spaces, we need to break down each sentence into a list of words through tokenization, while clearing up all the messy text in the process.

Creating Bigram and Trigram Models

Bigrams are two words frequently occurring together in the document. Trigrams are 3 words frequently occurring. Gensim's Phrases model can build and implement the bigrams, trigrams, quadgrams and more. The two important arguments to Phrases are min_count and threshold. The higher the values of these param, the harder it is for words to be combined to bigrams.

Some examples in our example are: 'health_insurance', 'health_care', 'white_house' etc.

Exploratory analysis

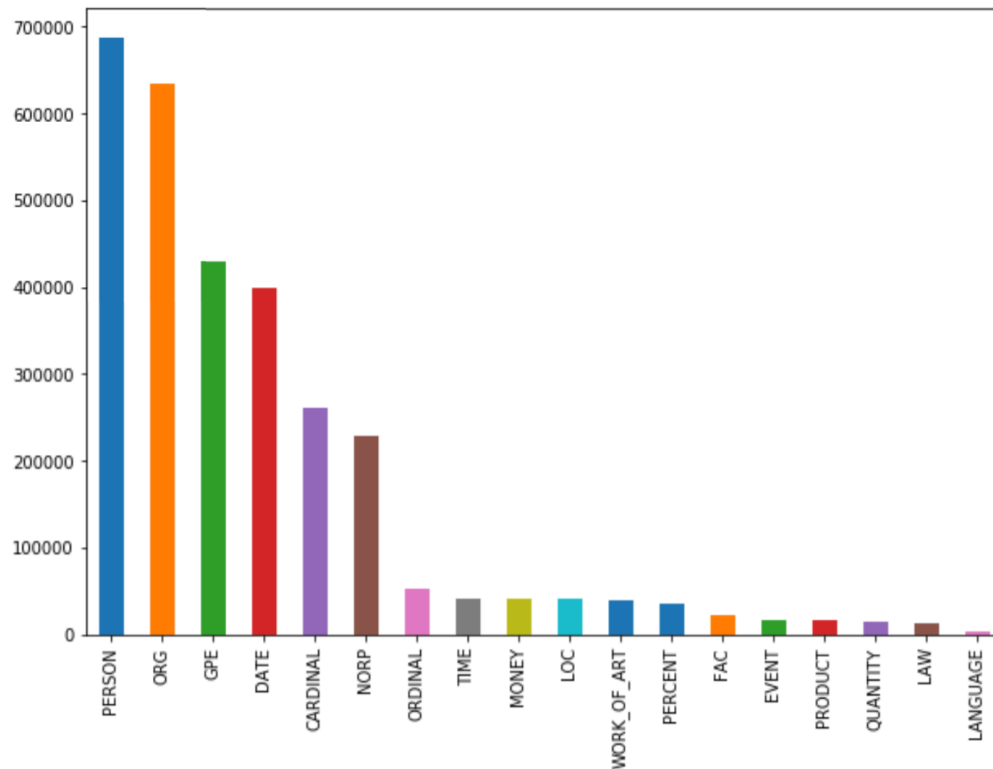
As of v2.0, spaCy supports models trained on more than one language, in this case is English. This is especially useful for named entity recognition.

To list a few examples: (Annotation Specifications <https://spacy.io/api/annotation#named-entities>)

TYPE	DESCRIPTION
PERSON	People, including fictional.
NORP	Nationalities or religious or political groups.
FAC	Buildings, airports, highways, bridges, etc.
ORG	Companies, agencies, institutions, etc.
GPE	Countries, cities, states.
LOC	Non-GPE locations, mountain ranges, bodies of water.
PRODUCT	Objects, vehicles, foods, etc. (Not services.)
EVENT	Named hurricanes, battles, wars, sports events, etc.

I first looked at the distribution of the number of the articles from different publications
And found out Breitbart, New York Post and NPR has the most number of news articles, so I
use them for the target of comparison.

After that I observed the distribution of types of predefined words by spaCy,
Which is listed below:



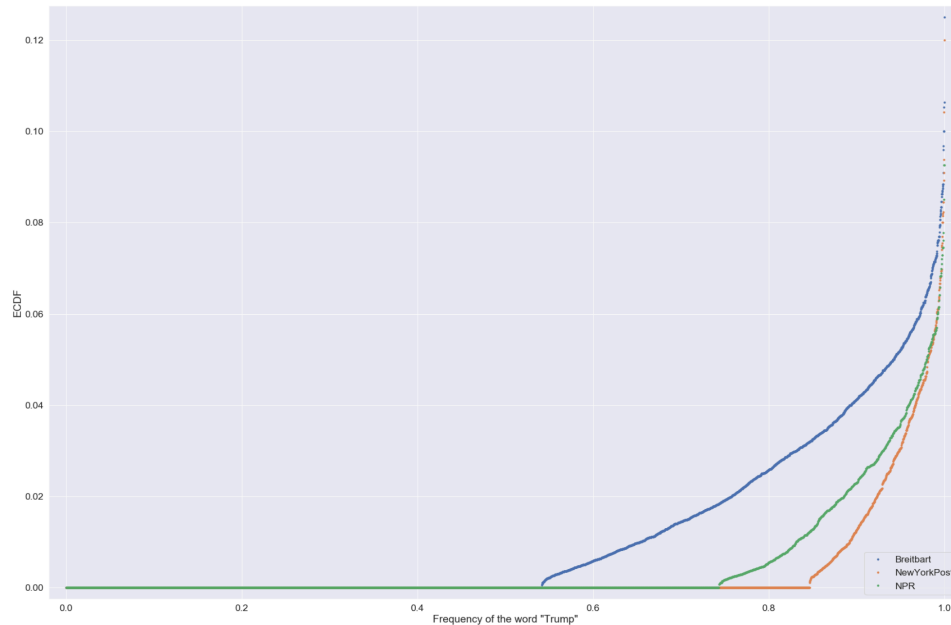
And I decided to look at some most occurring words in each media source.

From Breitbart: trump, state, people, president, news

From New York Post: year, time, people, day, trump

From NPR: people, year, trump, time, state

Since 'Trump' is in all three medias, but the frequencies of occurrence are different, I want to observe if such difference is statistically significant to notice such fact.



Statistical analysis

We can see that the differences of the occurring frequency of the word 'trump' from three different publications are much clearer in the ECDF. Breitbart has almost half of the articles contains the word 'trump'.

New York Post and NPR has more similar frequency for the word 'trump', if we can conclude that the frequency difference between them are statistically significant, then we can also conclude the difference between these two (New York Post, NPR) and Breitbart are much more statistically significant given its wider margin.

Preform hypothesis test

$H_0: p = 0$ $H_1: p \neq 0$

Null hypothesis:

There is no difference between the occurring frequency of the word 'trump' between New York Post and NPR (frequency difference is 0, $p = 0$)

Alternate hypothesis:

There is difference between the occurring frequency of the word 'trump' between New York Post and NPR (frequency difference is not 0, $p \neq 0$)

```
difference of means = 0.0016717863221870416  
p = 0.0
```

We get a p-value of 0.0, which suggests that there is a statistically significant difference. We got a difference of 0.00167 between the means. You should combine this with the statistical significance. The difference of 0.00167 in frequency is substantial.