

PanTop: Pandemic Topic Detection and Monitoring system

Yangxiao Bai, Kwanghee Won, Kaiqun Fu

Department of Electrical Engineering and Computer Science

South Dakota State University

Brookings SD, USA

{bai.yangxiao,kwanghee.won,kaiqun.fu}@sdstate.edu

Abstract—Diverse efforts to combat the COVID-19 pandemic have continued throughout the past two years. Governments have announced plans for unprecedentedly rapid vaccine development, quarantine measures, and economic revitalization. They contribute to a more effective pandemic response by determining the precise opinions of individuals regarding these mitigation measures. In this paper, we propose a deep learning-based topic monitoring and storyline extraction system for COVID-19 that is capable of analyzing public sentiment and pandemic trends. The proposed method is able to retrieve Twitter data related to COVID-19 and conduct spatiotemporal analysis. In addition to spatiotemporal data mining, the system implements an attention analysis on Twitter users. Furthermore, a deep learning component of the system provides monitoring and modeling capabilities for topics based on the most advanced natural language processing models. The proposed system includes a user-interactive visualization component that provides audience monitoring and analytical toolboxes. Case studies found in abundance in our proposed system can justify empirical analysis. Ultimately, we believe that our proposed system accurately reflects how public sentiments change over time along with pandemic topics.

Index Terms—social media analysis, topic modeling, storyline generation

I. INTRODUCTION

The rampaging COVID-19 pandemic has greatly affected emotion and everyone’s daily life in the past two years. The general public’s emotions went through a roller coaster-like up-and-downs: shocked by the deadly known disease, overwhelmed by misinformation, upset by the mitigation policies, and regaining confidence with the development of vaccines. Ubiquitous user-input content on social media and online services have generated a tremendous amount of information that can reflect the users’ emotion and opinion towards social events. With the abundance of the generated social media data during the pandemic, more users intend to express their emotions and opinions about COVID-19-related social events, such as the mitigation policies or the invention of the vaccines. Such dramatic social media data increase provides great research opportunities in social media mining and natural language processing.

Topic modeling analysis is one of the target research areas in social media analysis for COVID-19. In the case of pandemics, topic modeling analysis provides a reference for the development of epidemic mitigation measures and evaluation techniques. In the context of sentiment analysis,

some previous studies accomplished success by using the machine learning technique. For instance, previous work [?] used a deep learning-based sentiment analyzer to collect and curate the COVID-19 dataset. Attention mechanisms were applied to characterize the word-level interactions within a local and global context to capture the semantic meaning of words. Among other logical approaches, Chakraborty et al. [?] propose the implementation of fuzzy logic for taming the fuzziness of sentiments. However, the existing works in social media sentiment analysis lack the consideration of spatial and temporal factors. Our proposed system provides both spatial and temporal aspects of sentiment analysis under the topic of COVID-19.

Topic detection and modeling are the second focus of our proposed system. Topic detection is an important part of addressing the said problem to separate points of interest among many candidates. Several noble studies have been conducted in the last few decades to mine the most appropriate topic associated with a text phrase. A topic graph-based approach was proposed by Batool et al. [?] where a topic graph from vectorized Twitter data was proposed with term frequency as a heuristic and social relationship between the virtual users. Other work [?] proposed a time-dependent burst detection technique that focuses on two and three-word data acceleration as a spread tendency to make early detection based on the keywords. Abd-Alrazaq et al. [?] used Latent Dirichlet Allocation (LDA) for topic modeling on English twitter data. Gencoglu et al. [?] used a large Twitter dataset to analyze semantic topic clusters using Language Agnostic BERT Sentence Embeddings (LaBSE). We combine the advantages of LDA and neural network models to develop a transformative model for topic discovery and monitoring components in the system.

The storyline generation problem was first studied by Kumar et al. [?] as a generic redescription mining technique, by which a series of redescription between the given disjoint and dissimilar object sets and corresponding subsets are discovered. Storytelling is an efficient way to solve the issue of information overload. By extracting critical and connected entities, the original document is structurally summarized. Current works contain two categories: Textual Storytelling [?], [?], [?], [?], [?], [?] and Visual Storytelling [?], [?], [?]. A storyline generation component is deployed in the proposed

system. The storyline generation is capable of summarizing the highlighted COVID-19-related events with social media data, and a visualization component is also developed to demonstrate the comparisons between the official releases of the mitigation events and the identified topic events from social media.

In this paper, we present the Pandemic Topic Detection and Monitoring system (PanTop), which is capable of 1) collecting and retrieving social media data on COVID-19-related topics, 2) detecting hidden trends of topics from social media posts, and 3) generating and visualizing storylines for the extracted COVID-19-related topics. This paper is structured as follows: Section II discusses the architecture of the PanTop system; Section III illustrates the experiment and data of the system; Section IV demonstrates the case studies revealed by the PanTop system; in section V, we conclude our discoveries from the proposed PanTop system.

II. METHOD AND DATA

From a holistic view of research method (Figure 1), it consists several relational parts of the dataset import, dataset Processing, and machine model training for topic clustering

In dataset import process, since the dataset is un-hydrated, we needed to “hydrate” the data by building a process that will fetch the tweet content by querying with the Tweet ID. After hydrating, the data are saved in a local storage, and by using Tweets API the data are completed to be analyzed and stored in the Elasticsearch server. In dataset processing step,

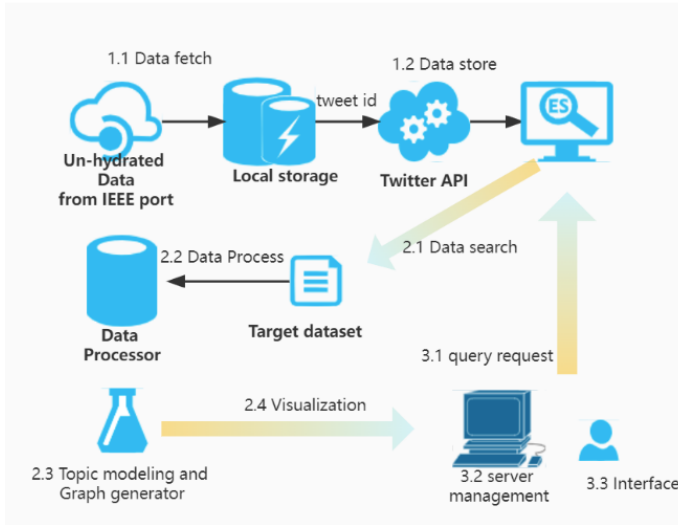


Fig. 1. System structure

after importing data into Elasticsearch server in batches, we choose a research scope and we get the target dataset from Elasticsearch. By doing spatial query and filtering task, the target dataset can be visualized into a map. Furthermore, after processing the data by using techniques such as the regex filtering, tokenization and tagging, and removing stop words, it is possible to apply a topic clustering model to get topics.

A. Sentence vector generation

We propose to use a LDA-BERT based method to generate vector for each tweet. First, we break down tweets into body and hashtags. For the body part, we tokenize the sentences and train them with BERT. Then we get vectors which contains hundreds of dimensions. For the hashtags part, we apply LDA on them and generate vectors of corresponding dimensions according to the number of topics to be divided into. To combine the information from two resource, we use autoencoder to compress combined vector.

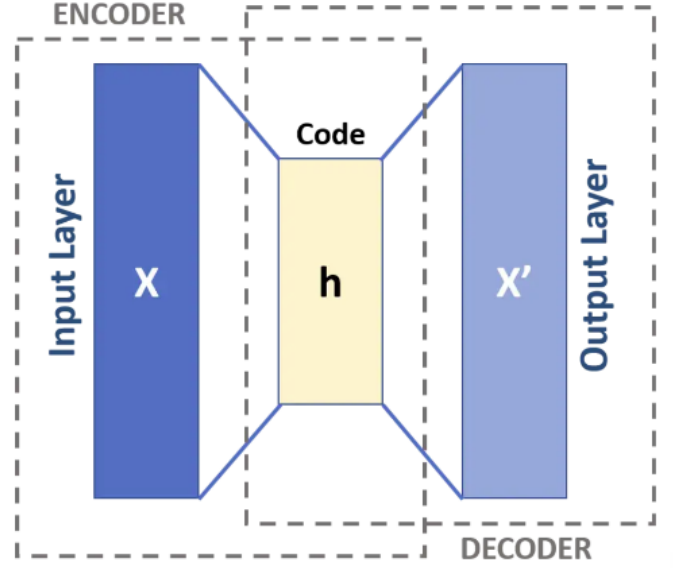


Fig. 2. structure of autoencoder

By this method, we succeed in reducing the dimensionality of the resulting vector while preserving topic features.

B. Topic Clustering

We use Time Series KMeans to do the clustering. Because each tweet is attached with publish time, we need to consider its temporal correlation while clustering.

C. Visualization

1) *Wordcloud:* Wordcloud is a graph to show the high frequency words. After clustering, we get labels which indicate which cluster the current tweet belongs to. Then we get all word lists under the same cluster and concat them. Finally the wordcloud will show higher frequency words in bigger font.

2) *Centroid sentence:* To study topic relevance, we calculate the cosine similarity between every tweet and cluster centroid point and pick up the nearest tweet as the most representative tweet.

3) *Storyline:* To generate a graph showing all topic relationships, we calculate the similarity between all the topic pair. Then sort the calculation results from high to low. Based on it, we set a threshold of edge. An edge will be generated between all topic pairs whose similarity is greater than the

threshold. Finally we import the point and edge information into Graphviz and get a graph.(Fig 5)

III. EXPERIMENTS

We apply our research method experimentally on a large spatial-temporal dataset collected from IEEE dataport. By comparing the result with real events in Coronanet which records government responses to coronavirus, we evaluate the plausibility of the method. The result shows several points of interest in a specific time period and the details within each topic.

A. Dataset Description and Experiment Setup

Dataset:

The dataset we used is CORONAVIRUS (COVID-19) GEO-TAGGED TWEETS DATASET form IEEE DataPort. This dataset includes English tweets related to the epidemic from 2020 to 2022 around the world. Due to the tweets spreading policy. We can't access the completed tweets content directly. The original data contains tweets ID and sentiment score calculated by dataset publisher such as:

ID	SCORE
12407280659839XXXXX	0.275

We use tweets ID to request twitter API in batch and get completed information, including content, UTC time, geographic location. All of these information are fully imported into the server after splicing with sentiment score. In this process, we will transfer the longitude and latitude to the geometry, then we can apply some spatial operator to handle them.

We use a transfer program named Twarc to batch fetch tweets in the dataset. Twarc is a python package that used to export tweets automatically. The main principle is that Twarc will make use of the registration information from twitter developer platform and get the permission from Twitter. Then Twitter can monitor the whole process of getting tweets. This way is completely legal and doesn't violate the privacy protocols set by Twitter. But it also cause some issues. If the tweets account is banned or the permission has been changed, we can't access the content anymore. Normally Twitter doesn't send any warning message but change the content of tweets to a prompt message. So some filtering process is necessary to ensure that the data set is not contaminated with irrelevant information.

We collect data from March 2020 to March 2022. and store all of these data to Elasticsearch with some batch programs, in order to facilitate subsequent steps to call on demand.

Scope of research: COVID-19 related tweets contain many types of topic. To avoid noise, we narrow our research from multiple perspectives. According to some statistical conclusions of the epidemic, we choose the time period May 2020 to August 2021 when the epidemic is most concerned. Geographically, we choose the continental United States as the research target. Besides, we set a keywords list to filter tweets, so as to limit the influence of irrelevant topics as much as possible.

(e.g. Posts about the city Corona instead of Coronavirus) All these search criterias are combined into one search request. The resulting sub-dataset will be used as the search object.

Preprocess:

The pre-processing process includes regex filtering, tokenization and tagging, removal of stop words. Regex filtering process is used to remove contents from which it is hard to obtain aluable information, such as URL, emoji, username, misspelled words. We also extract all hashtags from it.

Tokenization and tag are important to our project. We use nltk package to process the sentence in the tweets and decompose the content into single words. Then we can filter for parts of speech and get more reasonable results. The

```
[('i', 'NN'), ('m', 'VBP'), ('at', 'IN'), ('superior', 'JJ'), ('court', 'NN'), ('corona', 'NN'), ('branch', 'NN'), ('in', 'IN'), ('corona', 'NN'), ('', ''), ('ca', 'MD')]

[('joyeux', 'NN'), ('anniversary', 'JJ'), ('-', ':'), ('happy', 'JJ'), ('birthday', 'NN'), ('', ''), ('@', 'NNP'), ('titaylea', 'NN'), ('have', 'VBP'), ('fun', 'VBN'), ('and', 'CC'), ('enjoy', 'VB'), ('your', 'PRP$'), ('day', 'NN'), ('!', '!'), ('we', 'PRP'), ('love', 'VBP'), ('you', 'PRP'), ('.', '.'), ('@', 'VB'), ('sebsalcedo', 'JJ'), ('#', '#'), ('birthday', 'JJ'), ('#', '#'), ('bday', 'JJ'), ('#', '#'), ('joyeuxanniversaire', 'NN'), ('#', '#'), ('happy birthday', 'JJ'), ('#', '#'), ('friends', 'NNS'), ('#', '#'), ('bffs', 'NNS'), ('#', '#'), ('friendship', 'NN'), ('@', 'NNP'), ('corona', 'NN'), ('', ''), ('california', 'NN')]
```

Fig. 3. Tokenization and tag

purpose of removal of stopwords is to avoid distractions from common words on the topic. The nltk package provide a library of stopwords and we can set our own list of common words based on the theme. It's an important step to get good performance of machine learning models and we can adjust it to fit the current strategy we used.

Metrics:

We use silhouette factor to evaluate the performance of clustering. For storyline generation, we based on rationality of the time series and topic keywords.

Topic mining:

To evaluate the quality of generated topics, we focus on both the content of the topic itself and the distribution of the topic.

In the experiment, each topic consists of a series of related sentence vectors. We use high-frequency words in WordCloud to show the content of the topic. Besides, We also output one of the most representative sentences in each topic as supporting evidence.

The distribution of clusters in the vector space is another indicator of the quality of the results. The projection of vectors in two-dimensional space can intuitively see the distribution of the cluster. In order to verify the temporal relevance of the topic we also use the time label of the topic vector as the third axis to generate a three-dimensional vector distribution map.

Comparison Models:

We propose a combination of lda and bert methods to produce high-quality sentence vectors. To verify the advantages of this approach, we output the clustering results using bert and LDA separately. We make a comprehensive assessment based on the quality of the topics and timelines generated by each method

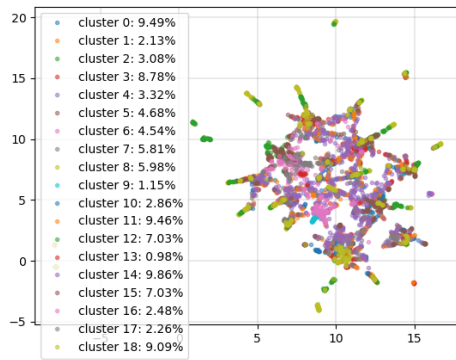


Fig. 6. 2d derivation for LDA

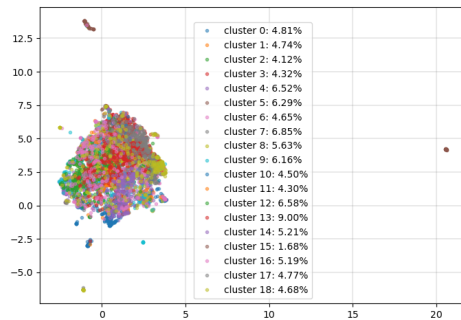


Fig. 7. 2d derivation for BERT

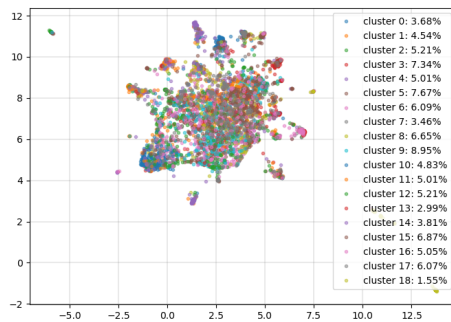


Fig. 8. 2d derivation for LDA-BERT

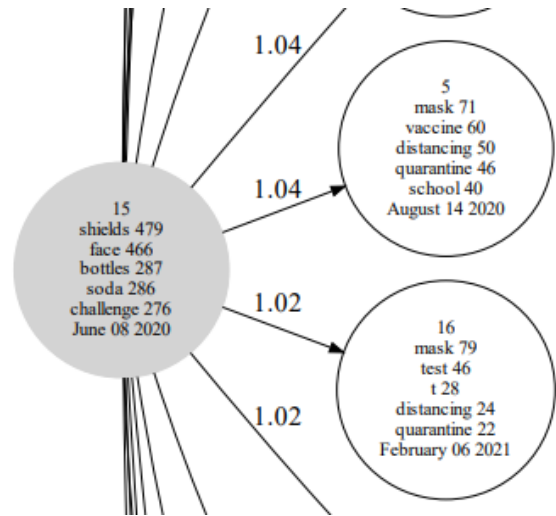


Fig. 9. slice of storyline for lda

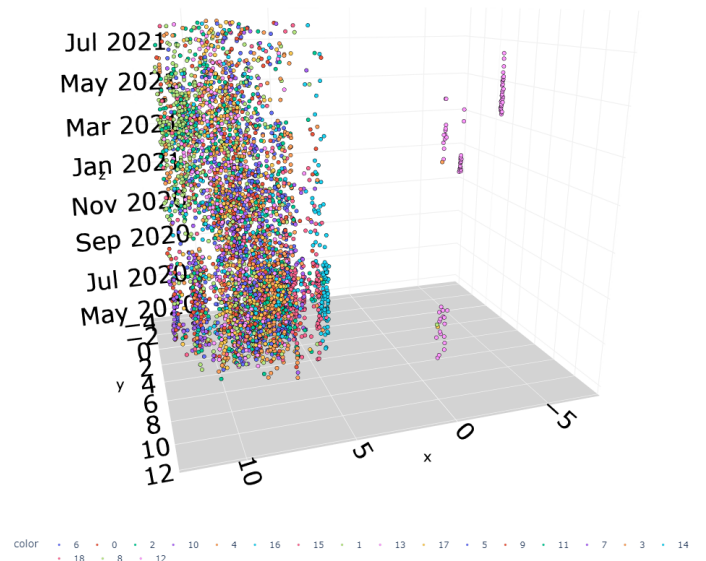


Fig. 10. 3D distribution graph

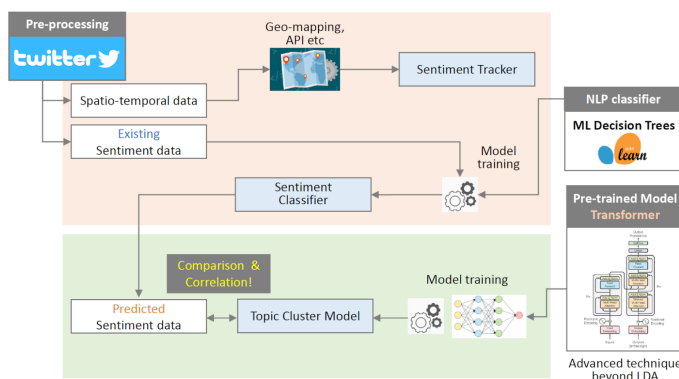


Fig. 11. Future plan

However it still has a long way to go in that it is too simple for disclosing some deeply meaningful implications. At present, since our model just provides the spatio-temporal view of the public opinion about an issue, it lacks further analysis about what specific topics are to be drawn from this. To dig into this issue, we will do detailed topic clustering task to classify how many different sub topics can be drawn from the overall topic of Covid-19 outbreak. Our topic for twitter dataset are filtered by the topic pandemic. However this is too simple. Therefore in the longer run, we should proceed to sub-topic clustering to get detailed areas under the single main topic. For example, they might be either government's specific policy about the pandemic or public opinion about

some vaccine company's business guideline, and etc.

In this project, we have chosen a static method for topic clustering which is Latent Dirichlet Allocation (LDA). Although this is static and an old method compared to the modern deep learning approaches, our result shows its relevancy in this domain is still present. The current research that are going on, a lot of the authors are using LDA to get a primary insight in the dataset that they compare with their model as a baseline. In the coming days, our research will focus on the modern approaches using deep learning to make compact and deployable topic clustering model.

One promising candidate for the model is the Transformer. This is a greatly applied model architecture in Natural Language Processing now, which is based on the core technology called the *attention mechanism*. Our team believe this should be an appropriate alternative current LDA to do the detailed topic clustering task (Figure ??).

V. CONCLUSION

In this project, our team carried out a machine learning-based spatial modeling of sentiment towards the pandemic. A rationale for this research comes from the increased interest on analyzing a lot of issues related to pandemic by using machine learning in academia. We have been pursued to combine the ML and Deep Neural Network (DNN) with our spatial-data mining technology so that this study contributes to the field by producing deeper implications.

As a dataset for sentiment analysis the importance of Twitter dataset can not be overstated. In recent days of DNN era, the public opinion expressed in twitter became a valuable resource for numerous tasks including such sentiment prediction. We paid attention on the applicability of DNN on the Twitter dataset for the pandemic issue. By using some applications such as Twarc and Elasticsearch we collected and stored such related dataset into our own repository, and spatial datamining task such as querying is followed. By doing such spatial query and filtering task, the target dataset could be visualized into a map. Furthermore, we develop a sentiment classifier by adopting the decision tree in Machine Learning. This module analyzes people's opinion in tweeter data and classifies it into positive or negative, by getting trained with pre-built sentiment score.

However, there are still another important task remained to achieve our team's final goal. Since our current model simply mine the spatio-temporal dataset and show how people's opinion is distributed differently by regions and time, it requires to do further research to draw deeper implications for contributing academia. Therefore, we plan to build a sub-topic clustering module using Transformer and compare its result with our predicted sentiment score. With adopting such a new NLP model architecture, it is expected to track a distribution of public opinion on detailed issues, e.g., government policy on social distancing, vaccine supply, etc., so that the study can have practical implication to decision makers.