

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/355226397>

Topic Modeling Using LDA and BERT Techniques: Teknofest Example

Conference Paper · September 2021

DOI: 10.1109/UBMK52708.2021.9558988

CITATIONS

2

READS

501

3 authors, including:



Ercan Atagün

Duzce University

15 PUBLICATIONS 22 CITATIONS

[SEE PROFILE](#)



Ahmet Albayrak

Karadeniz Technical University

40 PUBLICATIONS 96 CITATIONS

[SEE PROFILE](#)

Some of the authors of this publication are also working on these related projects:



KBU-15-DR-015 [View project](#)

Topic Modeling Using LDA and BERT Techniques: Teknofest Example

Ercan Atagün
Duzce University
Computer Engineering
Duzce, Turkey
ercan-atagun@hotmail.com

Bengisu Hartoka
Duzce University
Faculty of Technology
Duzce, Turkey
bengisu.hartoka@gmail.com

Ahmet Albayrak
Duzce University
Computer Engineering
Duzce, Turkey
ahmetalbayrak@duzce.edu.tr

Abstract—This paper is a natural language processing study and includes models used in natural language processing. In this paper, topic modeling, which is one of the sub-fields of natural language processing, has been studied. In order to make topic modeling, the data set was obtained by using the data scraping method, which has been very popular in recent years, over social media. The dataset is related to Teknofest competitions. The dataset was created by utilizing the Selenium library, one of the popular libraries used for the data scraping method. In order to be able to analyze on the prepared data set and to ensure the consistency of the clustering process, the text to be used before the analysis was preprocessed. After text preprocessing, clustering was performed on the data set with natural language processing techniques such as BERT and LDA.

Keywords—Data Scraping, LDA, Bert, Topic Modeling, Teknofest.

I. INTRODUCTION

Teknofest (Aviation, Space and Technology Festival) is a technology, aviation and space technology festival held in Turkey. It is the first and only festival organized by the Turkish Technology Team Foundation (also known as T3) and the Ministry of Industry and Technology, aiming to develop Turkey's national technologies. In general, it aims to promote and develop national technologies and to raise awareness of people in this field. While seminars, award-winning technology competitions, domestic technology exhibitions and an international entrepreneurship summit are taking place in the festival; Parachuting, SOLO TURK, Turkish Stars and similar local-international organizations air shows are also held. The festival is held one year in Istanbul and one year in one of the Anatolian provinces. In the festival, where many ministries and presidency are stakeholders, the stakeholder institutions vary according to the years, but these stakeholders are generally ASELSAN, ROKETSAN, Baykar, BMC, HAVELSAN, TÜBİTAK, STM, Turkish Aerospace(TUSAŞ), Turkish Airlines, Turkish Aeronautical Association, Türksat. and as supporters and participants such as Turkcell; is in establishment. In addition, the leading universities of our country provide support by providing academic consultancy services[1].

BERT is a complex trained model released by Google in 2018 that represents sentences in accordance with their word meanings. Since it can represent the sentence very well, it has performed very well in almost all sub-fields of NLP. Context-free models such as Word2vec and GloVe create a single word embedding for each word in the vocabulary [2]. Instead, contextual models such as BERT produce a representation of each word based on other words in the sentence. BERT captures these relationships bidirectionally as a contextual

model. BERT models have succeeded in gaining superiority over other models in many areas in the field of natural language processing in recent years [3]. These models are very complex in terms of their inner workings and require a significant amount of data and processing power for training. However, there are many ready-made BERT models. It is possible to achieve good results by using them and even by feeding them with additional training sets depending on the task to be done. The main goal of the BERT model is to create a fixed-length vector representing sentences, taking into account the meaning and order of words. BERT models differ depending on the language spoken. For Turkish, a German researcher “Stefan Schweter” published a Turkish-specific model, which can be accessed in Python with the transformers package [4], [5].

LDA is a probability-based topic modeling method. The model generates topics based on word weight from a set of documents. LDA is an example of topic modeling where each document is assumed to be a collection of topics and each word in the document corresponds to one of these topics. LDA is an unsupervised learning algorithm, it does not need predefined words. After the number of topics is determined, labels are assigned to the topics according to the classes[6]. Since the Bert model was used in this study, the use of GPU was needed. For this reason, Google Colaboratory was preferred as the coding environment for topic modeling. The most important reason for choosing Google Colaboratory is that it provides free GPU, unlike other cloud services. For topic modeling, vectors obtained from BERT and LDA were combined and clustering was applied [7], [8]. This study consists of data set cleaning and preparation, model development and experimental studies, and conclusions.

II. MATERIAL AND METHOD

A. Data Set

While preparing the data set used in this study, Visual Studio Code was used as the coding environment. Visual Studio Code is a fast, lightweight and software development tool that can be used on different operating systems. The tool, which has a simple text editor view, supports many programming languages such as Python, C/C++, C#, Javascript, thanks to plugins. In this study, the Python programming language was used [9].

The data set used in this study was prepared by using the Selenium library with the Web Scraping method, which has been very popular in recent years. Web Scraping is a technique used to retrieve data from web pages. It is an automated process in which the HTML codes of the web page are processed to extract data for manipulation, by executing on an

application's web page, copying it to a local database or a generated file for later analysis [9]. In this study, a bot was written to pull tweets on Twitter with Python programming language and HTML codes using Selenium library [10]. Using this bot, a total of 737 tweets (May 15, 2016 - November 20, 2020) were pulled from the Twitter pages of Teknofest and T3 foundation chairman of the board of trustees, Selçuk Bayraktar. Then, the captured tweets were saved in the csv file created in a sequential and orderly manner. In Figure 1, the block diagram showing the process steps is given.

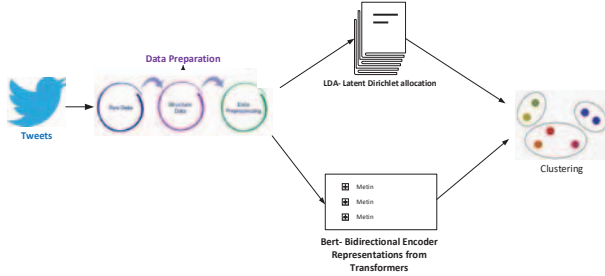


Fig. 1. Block diagram of the study

B. Text Preprocessing

In order to be able to analyze the text and to ensure the consistency of the clustering process, the text to be used before the analysis needs to be pre-processed. These processes are tokenization, removal of frequently used and meaningless stop words, conversion of all characters to lowercase, removal of information such as ip address or e-mail, numbers in tweets, but extracting them from the text, all kinds of punctuation and can be given as deletion of special characters. In Figure 1, these processes are given as data preparation. All these operations were carried out by defining separate functions. In this study, NLTK (Natural Language Toolkit) library was used for these operations. NLTK (Natural Language Toolkit) stands for natural language toolkit [11]. It is an open source natural language processing library [12]–[14] developed with the Python programming language to work with human language data and built with over 50 corpus and lexical resources under development. In Table 1, the raw text data prepared for this study is given before the text preprocessing process is applied.

TABLE I. ALGORITHM AND SCORE

<i>Id</i>	<i>Raw Data (tweet)</i>	<i>Text Preprocessing (tweet)</i>	<i>Tokenize (tweet)</i>
445	İHA dediğin doğa dostu olmalı... #MilliTeknolojiHamlesi	iha dediğin doğa dostu olmalı milli teknoloji hamlesi	'İha', 'dediğin', 'doğa', 'dostu', 'milli', 'teknoloji', 'hamlesi'
306	Milli Eğitim Bakanlığı'nın en eğlenceli okulu #TEKNOFEST'te @tcmeb alanında! #MilliTeknolojiHamlesi 20 Eylül Atatürk Havalimanı	milli eğitim bakanlığının eğlenceli okulu teknofestte alanında milli teknoloji hamlesi 20 eylül atatürk havalimanı	'milli', 'eğitim', 'bakanlığının', 'okulu', 'teknofestte', 'alanında', 'milli', 'teknoloji', 'hamlesi', 'eylül', 'atatürk', 'havalimanı'
38	Çevre ve Enerji Teknolojileri Yarışması #TEKNOFEST yarışmacıları temiz bir çevre, verimli	çevre ve enerji teknolojileri yarışması teknofest yarışmacıları	'enerji', 'teknolojileri', 'yarışması', 'yarışmacılar', 'i', 'çevre', 'verimli', 'gelecek', 'i'

<i>Id</i>	<i>Raw Data (tweet)</i>	<i>Text Preprocessing (tweet)</i>	<i>Tokenize (tweet)</i>
	bir gelecek için yarışıyorlar. @SankoholdingAs 24-27 Eylül Gaziantep Ortadoğu Fuar Merkez	temiz bir çevre verimli bir gelecek için yarışıyorlar 2427 eylül gaziantep ortadoğu fuar merkezi	'yarışıyorlar', 'gaziantep', 'fuar', 'merkezi'

Text preprocessing was applied in two stages. First, simple preprocessing was applied, all letters were converted to lowercase, non-alphanumeric characters were cleared by writing a function. Table 1 shows tweets with simple text preprocessing. As the last step of the text preprocessing applied in this study, the stop words of each tweet and the numerical values they contain were cleared and separated word by word, that is, tokenized. The result of this process is given in Table 1. After all these processes, the data set is ready to be given to the prepared model.

C. Model Development and Experimental Studies

This paper is a natural language processing study and has been studied on topic modeling, which is one of the sub-fields of natural language processing. There are generally two ways to define a topic. Similar structures in vector space can be clustered by resorting to hierarchical Bayesian models such as LDA (Latent Dirichlet Allocation) or by embedding the document in vector spaces. LDA is known to not do very well at processing short texts when there is not much text to model. Therefore, this study requires a technique that places the entire content of the sentence and can then cluster similar topics [15]–[17]. Therefore, creating vector representations and clustering may be the solution. Methods such as TD-IF, BERT have been tried to obtain vector representations from texts. TF-IDF is obtained by multiplying the term frequency (TF) of terms with the inverse document frequency (IDF). Class information of terms is not used in this weighting. Therefore, it is an unsupervised weighting method [8], [18], [19].

To cluster embeddings from BERT and LDA models and to find topics in clusters, it is necessary to ensure that documents with similar topics are clustered together. WordCloud was used for the visualization process. WordCloud is a data visualization technique used to represent text data where the size of each word indicates its frequency or importance. Important textual data points can be highlighted using a word cloud. WordCloud is widely used to analyze data from social networking websites. The libraries required to create WordCloud in Python are matplotlib, pandas and wordcloud [20], [21].

In this study, the result sets created by using the TF-IDF model were not balanced and sufficient for this study. Because TF-IDF is also word bag-based (ignoring grammar and word order), it loses contextual information, so a good enough result could not be obtained with the TF-IDF method. As given in Figure 2, the vectors obtained with TF-IDF do not distribute well enough and the processing results are not good in clustering.

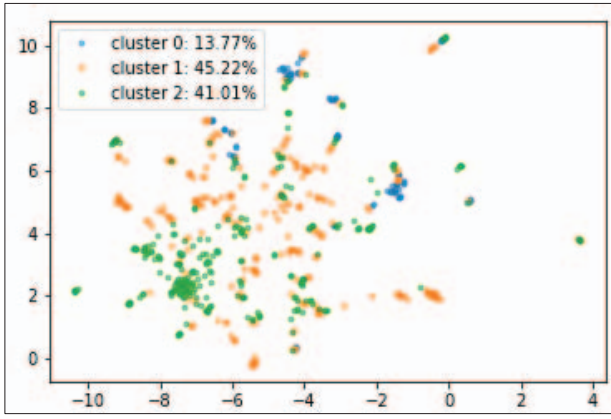


Fig. 2. Clustering result of vectors obtained with TF-IDF

The word distributions formed by the TF-IDF process are given in Figure 3 with word clouds.



Fig. 3. Word clouds obtained as a result of TF-IDF process

Vector representations of the tweets in the dataset were created using the BERT model. The Distilbert model was used for this study. This is because it provides a nice balance between speed and performance, and this package also includes a language model for Turkish [22]. When the vector representations from the BERT model were clustered, good enough results were not obtained. The clustering results made with the vectors obtained with the BERT model are given in Figure 4.

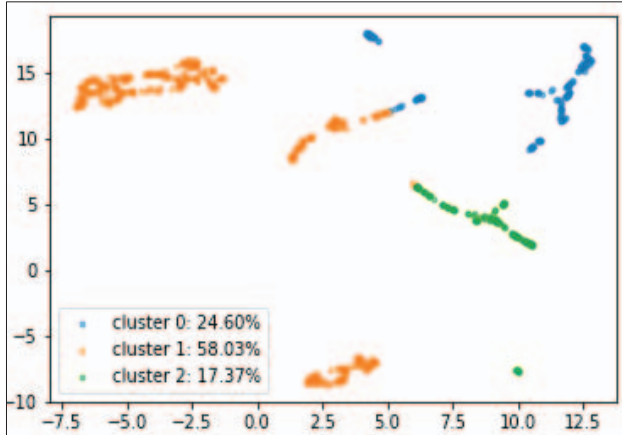


Fig. 4. Clustering result of vectors obtained with BERT

An important feature of the BERT model is that it places special start and end marks at the beginning and end of each text. The model outputs corresponding to the starting mark can be considered as the features that the model gives to the whole text [15]. As seen in Figure 4, the distribution obtained with the BERT model has been much more successful than the TF-IDF model, but it is still not a good enough result. In Figure 5, the word cloud obtained as a result of using the BERT model is given.



Fig. 5. Word clouds obtained as a result of the BERT model

In this study, the Clustering process was applied by combining the vector representations from BERT and LDA to make topic modeling, and the best result was obtained in this case. By combining LDA, BERT and clustering, semantic information can be preserved and contextual issue identification can be achieved. The model was designed by combining the probabilistic topic assignment vector of LDA obtained from the LDA model and the sentence vectors obtained from the BERT model [23]–[25]. Figure 6 shows the distribution resulting from Bert, LDA and clustering.

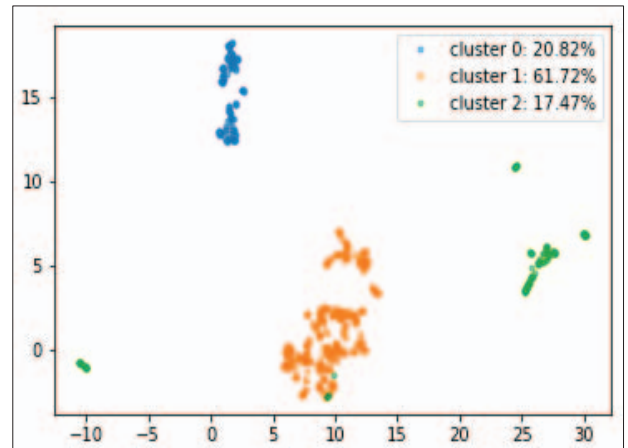


Fig. 6. Clustering result of vectors obtained with BERT and LDA

The word clouds of the topic modeling as a result of BERT, LDA and clustering are given in Figure 7.



- [24] D. Mazzei, F. Chiarello, and G. Fantoni, "Analyzing Social Robotics Research with Natural Language Processing Techniques," *Cognit. Comput.*, vol. 13, no. 2, pp. 308–321, 2021, doi: 10.1007/s12559-020-09799-1.
- [25] L. Lucy, D. Demszky, P. Bromley, and D. Jurafsky, "Content Analysis of Textbooks via Natural Language Processing: Findings on Gender, Race, and Ethnicity in Texas U.S. History Textbooks," *AERA Open*, vol. 6, no. 3, p. 233285842094031, 2020, doi: 10.1177/2332858420940312.