# Automatic Movie Rating by Using Twitter Sentiment Analysis and Monitoring Tool

Feriştah Dalkılıç
Department of Computer Engineering
Dokuz Eylül University
İzmir, Turkey
feristah@cs.deu.edu.tr
0000-0001-7528-5109

Ayşe Çam
Department of Computer Engineering
Dokuz Eylül University
İzmir, Turkey
ayse.cam@ceng.deu.edu.tr
0000-0002-7714-2117

*Abstract*— **Today, due to the intense use of social media platforms such as Twitter by all segments of today's technology, people have begun to share their views, ideas, and feelings through these media. It is possible to discover mighty valuable knowledge from this enormous resource. This study has emerged to assist users in making choices by evaluating emotions about TV series and movies that have recently appeared on social platforms, using ideas and feelings. The textual tweet data was preprocessed and cleaned of noise by using natural language processing techniques. Tweets were tagged using the Bert-based model according to the content of the Turkish TV series and movie comments, and their polarities were calculated. Machine learning models including Naïve Bayes (NB), Support Vector Machines (SVM), Random Forest (RF); Bagging and Voting, which are among the general ensemble algorithms, were trained for sentiment analysis by taking the obtained polarity values. The voting algorithm gives the best accuracy at 87%, while the Support Vector Machines give the best area under the receiver operating characteristics curve (AUC) of 0.96. A web application was developed by using Flask to monitor sentiment scores via hashtags (#).**

*Keywords—sentiment analysis, machine learning, natural language processing, social data science*

## I. INTRODUCTION

Today, due to the increasing internet use, users communicate through social networks. Social networks have become the easiest way people can express themselves and spread their ideas. People update their status in every activity they do, share their photos, and instantly share their thoughts about the TV series, movies and TV shows they watch. For this reason, issues such as the processing of data received from social networks gain importance day by day. Most of the evaluations in this study are on Twitter, new TV series and movies are on Netflix, AmazonPrime, BlueTV etc. published on platforms.

Several studies have been conducted to extract sentiment from social media data, especially from Twitter data in the last decades. But very few of these studies are in Turkish [1]. In this study, the tweets in Turkish that have hashtags (#) with movie titles were evaluated, classified, and rated as positive, negative, or neutral. Along with these, the total number of people who tweeted, the first user to use the hashtag showed how it affects women and men. Twitter API was used to access and collect Twitter data. The tweets were stored in MongoDb with NoSQL infrastructure. The following stages were applied to tweets pulled according to hashtags, respectively: Words were normalized and converted to lowercase. String tokenization was applied to texts, stop words were removed. Since the tweets were found without tags, they were tagged with certain models. A rating system was proposed to show the likes and dislikes rates over many different tweets and comments for new movies and series, such as IMDB ratings, on a single site.

## II. RELATED WORKS

In recent years, sentiment analysis and opinion mining studies have become a popular subject in the field of natural language processing (NLP). Social media reviews, forums, microblogs, and product reviews have been used in many of the studies carried out in this field. Twitter data has been analyzed to understand social media users' opinions on various issues in daily life. Zimbra et al. gave a detailed survey of Twitter sentiment analysis applications [2]. Wang et al. used Twitter data about the 2012 US presidential election to develop a system for real-time sentiment analysis [3]. Abalı et al. detected the problems of citizens and extracted the locations of complaints from Turkish tweets collected from the Aegean Region of Turkey [4]. Chakraborty et al. handled COVID-19 related tweets during two distinct pandemic times to analyze the sentiment tendency of tweets [5]. Pant et al. performed the prediction of Bitcoin's volatile price by analyzing sentiment on Twitter [6]. A sentiment analyses of the automotive industry tweets was presented by Shukri et al. to extract the polarity and emotions classification towards the automotive classes [7].

Several different techniques have been used in sentiment analysis and opinion mining studies hitherto. Some studies were conducted on lexicon-based methods while others were conducted on machine learning-based methods. These techniques have been combined in some recent studies.

Lexicon-based methods use a sentiment dictionary to detect the sentiment tendency of text. Quan and Ren conducted a feature-based sentiment analysis study on product reviews [8]. For feature extraction, they first determined all names in the interpretation as candidate traits, then they used Term Frequency-Inverse Document Frequency and Pointwise Mutual Information (PMI) methods together to determine the level of relationship of these candidate features with the product. After determining the product features, they used a commitment parser to detect emotion expressions and emotion polarities. They tested their proposed method on their comments on digital cameras, cell phones, mp3 players, and routers. In these tests, they achieved success rates between 61% and 89% in feature extraction and between 66% and 77% in emotion classification. In their study, Atan and Çınar used news texts published in different news sources in 2014 regarding BIST30 companies traded on Borsa Istanbul as a data set [9] and converted the expressions in the news content into numerical values with the help of a sentiment dictionary translated into Turkish. Then, the relationships between these

numerical scores and the company values formed in the market in the same period were analyzed. The main result that emerges is that there are significant relationships between the news published in financial markets and their emotional tone and financial values. Karagöz and Gürsoy used the tweets written about the programs broadcast on a TV channel in an eight-month period as a data set [10]. Sentiment analysis revealed that the messages about the channel and the program included in this data set contained emotions classified as positive, negative, or neutral. With this information, channel managers were able to make predictions for program managers and it was stated that they could develop relevant strategies in this direction.

In some of the recent studies, machine learning methods are used to identify the sentiment polarity in texts [11]. Both unsupervised learning-based clustering techniques and supervised learning-based classification techniques of machine-learning have been employed for sentiment analysis. Li et al. tried to determine the locations of incidents such as murder, accident, and disease by predicting the time and place using tweets. While ranking similar tweets, they used clustering methods. Finally, they used a Linear Regression (LR) model to score the importance of tweets [12]. Tuzcu used a data set containing reader comments taken from an online book sales site to make emotion classifications with classification algorithms [1]. An emotion classification was performed with Multi-Layer Perceptron (MLP), NB, SVM, and LR algorithms. Although all algorithms have done well, MLP was the algorithm that performs the best. Basiri et al. proposed a method based on a fusion of four deep learning and one traditional supervised machine learning model on coronavirus-related tweets for sentiment analysis [13]. They stated that their model is better than four individual deep models and the DistilBERT model.

Hybrid methods merge machine learning and lexicon-based methods to recognize text emotion. Torres developed a sentiment analysis service on student-generated comments in Spanish. The original system performs the natural language processing task to determine the sentiment associated with a text. SVM, Naïve Bayes, logistic regression, and decision tree machine learning models have been used for sentiment classification. The results showed that machine learning-based methods outperformed the original system [14]. Kaynar et al. used Naïve Bayes, Central Based Classifier, Multi-Layer Perception (MLP), and SVM for sentiment analysis of movie reviews in IMDB [15]. According to experiments, it is seen that MLP and SVM outperform both training and test sets. In the training dataset, MLP performed better with an 89.73% correct classification rate while both classifiers showed almost the same performance with a correct classification rate of around 75% for the test data set. Patel and Passi used the tweets about the 2014 World Cup soccer tournament held in Brazil to analyze the emotion of people all over the world by combining lexicon-based and machine learning methods [16]. They obtained the best accuracy of 88.17% by the Naïve Bayes classifier, and the best AUC of 0.97 by the random forest classifier.

## III. METHOD

In this study, we performed natural language processing techniques including word tokenization, stop word elimination, word stemming, and lemmatization. Along with these, spell checkers have been applied to correct spelling mistakes. While tagging tweets, the Bert model was activated with the support of Google Colab's GPU augmentation, and a training dataset of 3 sentiment moods was prepared. TF-IDF vectorizer was used to create the polarities of the words separated into their stems, and these digitized values were made ready for use in machine learning algorithms. The general flow of the proposed method is given in Fig. 1.
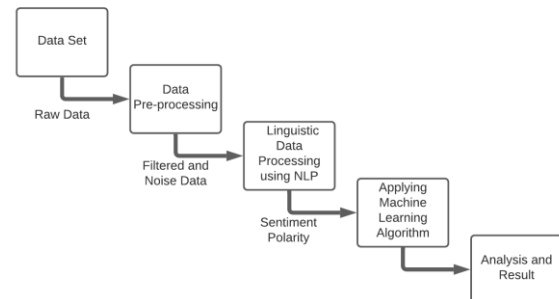


Fig. 1. Flowchart of the proposed method.

### A. Dataset

In the scope of this study, a dynamic data set consisting of tweets with hashtags related to movies and TV series that have been published recently was collected. Tweets were pulled according to keywords by accessing the Twitter API and using the Tweepy library. Datetime and base64 libraries were used to access instant tweets. MongoDB, a NoSQL database, was preferred to store large amounts of data and perform fast searches [17]. MongoDB keeps adding, updating, deleting, and searching non-relational data on the file system in Json format. Sentiment analysis studies were conducted with a total of 19,398 tweets.

### B. Preprocessing Tasks

The real meaning and sentiment of the text data can only be derived correctly if the data are cleaned from noise or irregular patterns. To perform analytical operations in tweets, words must be digitized and measurable. For this purpose, it is important to make the words as simple as possible to measure the frequency of the words.

TABLE I. ALGORITHM FOR THE PREPROCESSING OF TWITTER COMMENTS

| STEP | OPERATION |
|---|---|
| 1 | By using regular expression techniques, all URLs are replaced with the 'URL' keyword |
| 2 | All '@username' terms are replaced with the 'AT_USER' keyword |
| 3 | All #Hashtags and RT are eliminated from the text |
| 4 | Capitalization conversion to lowercase |
| 5 | Double spaces and double characters are corrected |
| 6 | Turkish characters [^A-Za-zığüşöç] are allowed |
| 7 | Predefined special characters (: n \| [ ] ; : {} − + ( ) < > ? ! @ # % *,) are eliminated |
| 8 | Comment texts are tokenized |
| 9 | Stop words are eliminated |

The crawled tweets contain URLs, hashtags \#", annotation \@" and retweets \RT" in addition to text data. Before applying machine learning techniques to the text data, text input must be tokenized. Preprocessing was performed by passing the data through the steps of capitalization conversion, elimination of punctuation marks, and elimination of stop words. The Turkish corpus of the NLTK library was used in the filtering of the stop words. The Word Punct Tokenizer module from the NLTK package was used to split the text into tokens and delete punctuation. Preprocessing steps applied to tweets are given in Table 1.

There are many natural language processing libraries developed for Turkish, the Zemberek project was preferred for its features and ease of use. After cleaning the data by using several well-known techniques (outlier detection, …), the data got ready to be labeled. There may be some spelling mistakes or different spelling in raw words due to the Twitter environment. Using SpellChecker from the Zemberek library, it was checked whether the words were spelled correctly. If a spelling mistake was detected, the first suggestion was used to correct it. It was also added to the temporary empty list.

The initial unconjugated form of words is called lemma. So, each word has to be transformed its lemma form before labelling. This is the essential aim of lemmatization that parses the lemma form of the word by finding out the suffixes and/or prefixes added or prepended to the word, and in our study Zemberek library was our valuable assistant.

### C. Data Labeling

After the data was pre-processed; 19,398 tweets were tagged as positive, negative, and neutral in order to train the model. At this stage, the BERT model developed by Google was used [18]. First, a small model was created using the BERT algorithm, and this model was run on 19,398 tweets. In contrast to the superficial language processing that goes from right to left and left to right, the BERT algorithm found the relationship of each word with the other word, resulting in complex but more accurate hashtag tweets.

After tagging the tweets, Fig. 2 is prepared by using the ratio of positive, negative and neutral results. Tweets that are only in the Turkish language were analyzed and the tweets in other languages were ignored. The reason why there are so many negative tweets is due to users' reactions to events in the scenes. Some sample pre-processed and tagged tweets can be seen in Fig. 3.
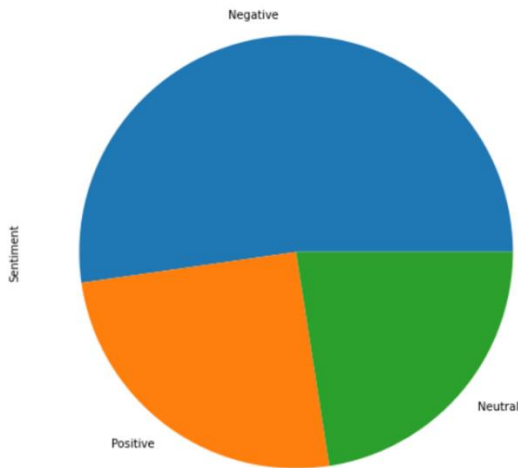
Fig. 2. Sentiment distribution of tweets.

| Tweet | Sentiment | Analysis |
|---|---|---|
| yalan söyleyeyim üzgünüm yanım kopmuş kanamış ... | -1 | Negative |
| ilker boşandıktan gerçekten zafer işareti yapt... | -1 | Negative |
| hande saf kötüsün ▓▓ kızı saçını başını yolas... | -1 | Negative |
| şte budur adaletsiz guen osmanbey çocuklarakıy... | 1 | Positive |
| ▓▓ kuzen işi var elanın odasında masumiyet ... | -1 | Negative |
| ... | ... | ... |
| beynim yandı sizce suçlu masumiyet | -1 | Negative |
| kisinde izlemiyorum | 0 | Neutral |
| linç kültürünü iyi anlatan dizi bence kız kend... | -1 | Negative |
| masumiyet hülya avşar filtreli görünce kime be... | 0 | Neutral |
| kadının büyük düşmanı yine kadın masumiyet | -1 | Negative |

Fig. 3. Sample pre-processed and tagged tweets

### D. TF-IDF Model

Before moving on to classification modeling with machine learning, data should be represented numerically. In this part of the study, we digitized the data using the term frequency-reverse document density (TF-IDF) model. The more the term is repeated in the document, the higher the TF-IDF value. In the coding part, TfidfVectorizer class was used to convert TF-IDF feature matrix. Bigrams and unigrams were included using ngram_range in TfidfVectorizer.

Equation (1) and (2) are used to calculate the total positive score (TotalPosScore) and total negative score (TotalNegScore), where $n$ is the number of terms, and $t$ denotes the tweet [16].

$$TotalPosScore_t = \sum_{s=1}^{n} TotalPosScore + PosScore_s \quad (1)$$

$$TotalNegScore_t = \sum_{s=1}^{n} TotalNegScore + NegScore_s \quad (2)$$

As given in Equation (3) sentiment polarity ($Polarity_{sa}(t)$) of a tweet can be a negative value, neutral (0) or a positive value [16]. The sentiment polarity is calculated for each of the tweets.

Fig. 4 demonstrates the frequency distribution of top 50 tokens from a typical download of 10,700 tweets. Positive terms such as "masumiyet" and "güzel" appear frequently in tweets. Series names and artist names are also frequently observed.

### E. Algorithms

The classifiers were trained using pre-labeled Twitter data to precisely label the moods associated with the text, thus achieving the highest possible accuracy. We split the dataset as training and test sets by the ratio of 80% and 20%. Some of the classifiers have a tendency to cause overfitting by learning the detail and noise in the training dataset and show worse performance on a new dataset. We used 10-fold cross-validation to prevent the overfitting of classifiers.

Naïve Bayes, SVM, Random Forest algorithms, Bagging and Voting ensemble learning techniques were selected as machine learning algorithms. While the "sklearn" library in Python was used for machine learning algorithms, the "zeyrek" and "jype" libraries were used for zemberek.

Sklearn.svm library allows us to change the cost kernel and gamma values for the SVM algorithm. In the sklearn.model_selection library, parameters such as max depth in the tree, estimators, and maximum feature were set to the best estimator in the RF algorithm and run with maximum efficiency.

### F. User Interface

In the scope of this study, a web-based user interface was designed for an easy and effective user experience. Users can display the general emotion distribution of current movies and TV series graphically and statistically on a dashboard.

Another page was designed to monitor the most popular hashtags. Users can see the top hashtags and how many tweets related to a specific hashtag have been posted. A hashtag search page was supplied for the users to get more information about the hashtags they were interested in. By using the user interface given in Fig. 5, users can get details of hashtags they searched. On this page, the user can access the sentiment analysis results of the tweets, the username of the users who post the tweets, and the number of followers. The user can also view statistically how many positive, negative, or neutral tweets are related to the hashtag.

$$Polarity_{sa}(t) = \begin{cases} Positive \text{ or } 1, if \text{ } TotalPosScore(t) > TotalNegScore(t) \\ Negative \text{ or } -1, if \text{ } TotalPosScore(t) < TotalNegScore(t) \\ Neutral \text{ or } 0, otherwise \end{cases} \quad (3)$$

TABLE II.  EXAMPLE OUTPUT OF SENTIMENT ANALYSIS

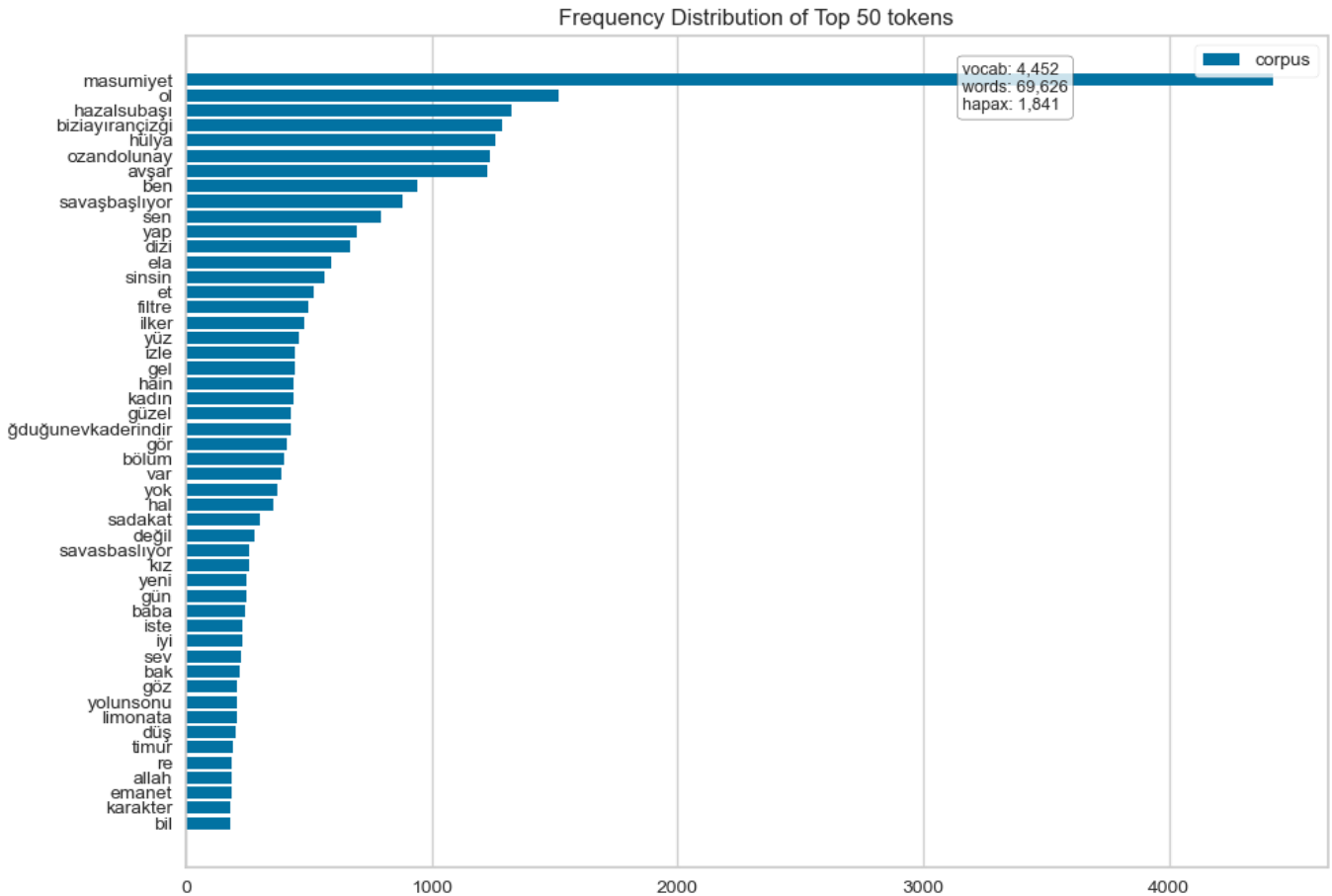| Tokens | Total PosScore | Total NegScore | Sentiment Polarity |
|---|---|---|---|
| ['Hülya', 'Avşar', 'a', 'seda', 'sayan', 'instagram', 'filtresi', 'yap', 'bakma'] | 0.082781 | 0.916643 | Negative |
| ['grup', 'kucaklaş', 'ilker', 'masumiyet', 'harika'] | 0.779853 | 0.050602 | Positive |
| ['mehdi', 'son', 'bölüm', 'allah', 'biz', 'sabır', 'ver', 'yolunsonu'] | 0.039630 | 0.532317 | Negative |



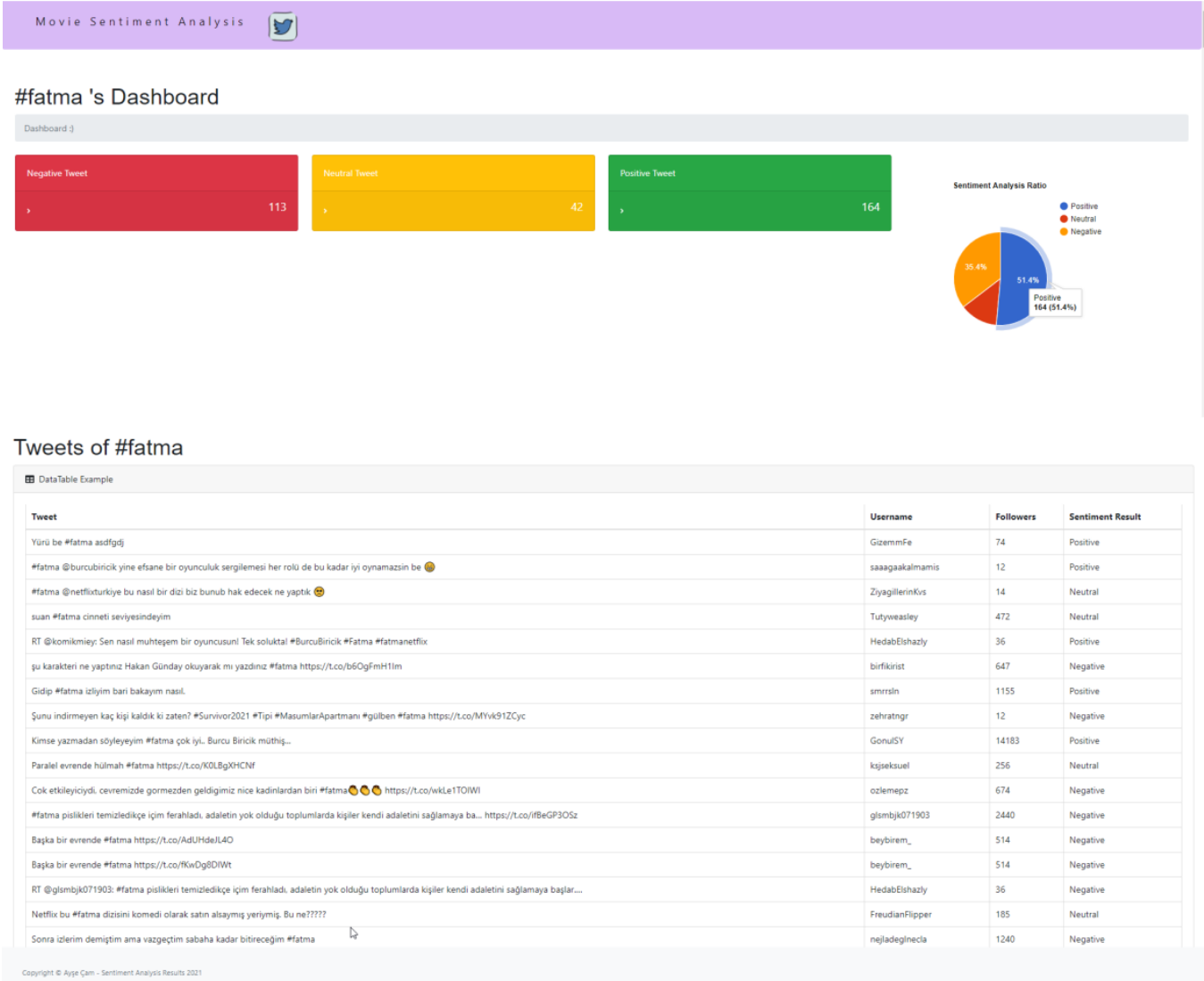Fig. 4.   Frequency distribution of top 50 tokens

Fig. 5. Tweet-based and overall sentiment results for the hashtag #fatma.

## IV. RESULTS

We labeled the tweets automatically by using BERT model and obtained 7,606 positive, 5,319 neutral, and 6,075 negative-labeled reviews. Machine learning algorithms have a tendency to ignore the minority class. We used an over-sampling method to avoid this problem. Naïve Bayes, SVM, random forest, and Voting and Bagging machine learning techniques are the techniques used to train the dataset and each technique creates slight differences in the result. These techniques are the well-known techniques used in sentiment analysis.

In the dataset, 80% of the data is used for training and the remaining (20%) for testing. Performance results of different techniques are given in Table 3 where 10-fold cross-validation technique is used to get realistic results. Parameter tuning was applied for each classifier and the best performance of each was obtained.

The Voting and SVM algorithms seem to give the best AUC values as 0.962 and 0.963, respectively. The best accuracy was obtained by the Voting algorithm as 87%. We obtained more accurate prediction results by combining lexicon-based and machine learning methods.

TABLE III. PERFORMANCE RESULTS FOR SENTIMENT ANALYSIS

|  | Precision | Recall | Accuracy | F-Measure | AUC |
|---|---|---|---|---|---|
| **SVM** | 0.88 | **0.86** | **86%** | **0.86** | **0.963** |
| **Naïve Bayes** | 0.86 | 0.82 | 85% | 0.83 | 0.950 |
| **Bagging** | 0.88 | 0.81 | 85% | 0.83 | 0.959 |
| **Voting** | **0.89** | 0.83 | 87% | **0.86** | 0.962 |
| **RF** | 0.87 | **0.86** | **86%** | **0.86** | 0.954 |

There are a few studies that deal with the sentiment prediction of text in the Turkish language. Atan and Çınar made dictionary-based emotion scoring and reached a correlation value of 0.79 between the emotional tones of the news and the market value of the companies [9]. In another study, Karagöz and Gürsoy, took the tweets written about the programs broadcast on a television channel over an eight-month period as a data set. By using semi-supervised technique, they created two separate dictionaries with positive and negative words and reached 68% accuracy by looking at the positive and negative excess with the help of excel.

The most obvious difference of this study from the existing studies on Turkish is the automatic labeling of tweets according to their sentiments. While the text data in other studies were labeled with the help of humans, the unlabeled data were first labeled with the Bert model, and then the success rate was increased up to 87% by using n-grams and using more than one machine learning algorithm in this study. When compared with the existing sentiment analysis studies performed on Turkish texts, our study has a promising performance.

## V. CONCLUSION

In this study, emotion analysis was performed on the linguistic tweet dataset that contains comments about Turkish TV series and movies. The first step after collecting the tweets was cleaning the noise and removing outliers. By pruning and shaping the data, the data became a smaller and well-been dataset. Sentiment scores were calculated using Word stemming, lemmatization, and n-gram techniques. Tweets were tagged automatically using the Bert-based model according to their content. TF-IDF feature matrix was constructed using bigrams and unigrams and textual data was prepared for machine learning.

Machine learning algorithms including Naïve Bayes, SVM, Random Forest, Voting, Bagging were used to classify tweets as positive, neutral, or negative sentiment, and the performance of classifiers was measured for the accuracy and reliability of the classification models. Although all the classifiers performed well on the prepared dataset, when we handle the AUC metrics, SVM was the best classifier with the 0.963 AUC. Voting was the other successful classifier with 87% accuracy.

Consequently, in this study, a model that made more accurate predictions with SVM and Voting algorithms was created, and the emotional states of the audience of TV series and movies were displayed by the developed user interface.

## REFERENCES

[1] S. Tuzcu, "Çevrimiçi Kullanıcı Yorumlarının Duygu Analizi ile Sınıflandırılması," *Eskişehir Türk Dünyası Uygulama ve Araştırma Merkezi Bilişim Dergisi*, 1(2), 1-5, 2020.

[2] D. Zimbra, A. Abbasi, D. Zeng, and H. Chen, "The state-of-the-art in Twitter sentiment analysis: A review and benchmark evaluation," *ACM Transactions on Management Information Systems (TMIS)* 9.2, 2018, 1-29.

[3] H. Wang, D. Can, A. Kazemzadeh, F. Bar, and S. Narayanan, "A system for real-time Twitter sentiment analysis of 2012 us presidential election cycle," *In Proceedings of the ACL 2012 system demonstrations*, pp. 115-120, July 2012.

[4] G. Abalı, E. Karaarslan, A. Hürriyetoğlu, and F. Dalkılıç, "Detecting citizen problems and their locations using twitter data,"*In Proceedings of the 6th International Istanbul Smart Grids and Cities Congress and Fair (ICSG)*, 2018, pp. 30-33, doi: 10.1109/SGCF.2018.8408936.

[5] K. Chakraborty, S. Bhatia, S. Bhattacharyya, J. Platos, R. Bag, and A. E. Hassanien, "Sentiment Analysis of COVID-19 tweets by Deep Learning Classifiers—A study to show how popularity is affecting accuracy in social media," *Applied Soft Computing*, 97, 106754, 2020.

[6] D. R. Pant, P. Neupane, A. Poudel, A. K. Pokhrel, and B. K. Lama, "Recurrent neural network based bitcoin price prediction by twitter sentiment analysis," *In Proceedings of the IEEE 3rd International Conference on Computing, Communication and Security (ICCCS)*, pp. 128-132, October 2018.

[7] S. E. Shukri, R. I. Yaghi, I. Aljarah, and H. Alsawalqah, "Twitter sentiment analysis: A case study in the automotive industry" *In Proceedings of the IEEE Jordan Conference on Applied Electrical Engineering and Computing Technologies (AEECT), IEEE*, pp. 1-5, November 2015.

[8] C. Quan and F. Ren, "Unsupervised product feature extraction for feature-oriented opinion determination," *Information Sciences*, 272, 16-28, 2014.

[9] S. Atan and Y. Çınar, "Borsa İstanbul'da finansal haberler ile piyasa değeri ilişkisinin metin madenciliği ve duygu (sentiment) analizi ile incelenmesi," *Ankara Üniversitesi SBF Dergisi*, 74.1, pp. 1-34, 2019.

[10] B. Karagöz and U. T. Gürsoy, "Adaptif Öğrenme Sözlüğü Temelli Duygu Analiz Algoritması Önerisi," *International Journal of Informatics Technologies*, 11(3), pp. 245-253, 2018.

[11] B. Pang, L. Lee, and S. Vaithyanathan, "Thumbs up? Sentiment classification using machine learning techniques," *arXiv preprint*, cs/0205070, 2002.

[12] R. Li, K. H. Lei, R. Khadiwala, and K. C. C. Chang, "Tedas: A twitter-based event detection and analysis system," *In Proceedings of the IEEE 28th International Conference on Data Engineering*, pp. 1273-1276, April 2012.

[13] M.E. Basiri, S. Nemati, M. Abdar, S. Asadi, and U.R. Acharrya, "A novel fusion-based deep learning model for sentiment analysis of COVID-19 tweets,"*Knowledge-Based Systems*, 228 (2021): 107242.

[14] M. J. D. Torres, "Contributions to Social Learning Analytics based on Sentiment Analysis of Students Interactions in Educational Environments," Doctoral dissertation, Universidad De Las Américas Puebla, 2019.

[15] O. Kaynar, Y. Görmez, M. Yıldız, and A. Albayrak, "Makine öğrenmesi yöntemleri ile duygu analizi," *In Proceedings of the International Artificial Intelligence and Data Processing Symposium*, pp. 17-18, 2016.

[16] R. Patel and K. Passi, "Sentiment analysis on Twitter data of world cup soccer tournament using machine learning," *IoT*, 1(2), pp. 218-239, 2020.

[17] MongoDB (2021, December 18), Retrieved from https://www.mongodb.com/.

[18] J. Devlin, M. W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," *arXiv preprint arXiv:1810.04805*, 2018.