

Sentiment Analysis of Arabic Tweets on the Great March of Return using Machine Learning

Saad Tareq Daher¹, Ashraf Yunis Maghari¹ and Hussam Fares Abushawish²

¹Islamic University of Gaza, Palestine, {sdaher@iugaza.edu.ps, amaghari@iugaza.edu.ps}

²Digital Palestine Technical College, Palestine, {hussam.as@gmail.com}

ABSTRACT

Social media platforms such as Twitter and Facebook are becoming powerful sources of people's perception of major events. Most people use social media to express their views on various issues and events and develop their information on a diverse economic, political, technical, social and occurrences related to their life. The overarching aim of this paper is to apply machine learning techniques to extract Arab users' opinions from 500 Arabic tweets on the Great March of Return rallies in the Gaza strip (Gaza border protests) collected over a two years span from 2018 to 2019. The majority of Sentiment Analysis (SA) studies concentrate on the English language, while other popular languages, such as Arabic, are seldom covered. In addition, on the Internet, publicly accessible Arabic datasets are hardly found. Three Arabic sentiment analysis datasets were used to train and evaluate four machine learning algorithms, namely, Support Vector Machine, Logistic Regression, Decision Tree, and Neural Network. In term of accuracy, logistic regression outperformed the other three algorithms with a percentage of 83%. Application of logistic regression on the sample tweets revealed that 85.8% of the tweets opposed the Great March of Return, whereas 14.2% of the tweets supported it.

Keywords: Opinion mining, sentiment analysis, Arabic language, machine learning.

I INTRODUCTION

The Great March of Return (GMR) rallies (Arabic: مسيرة العودة الكبرى) started on the 30th of March, 2018 with a series of demonstrations held on every Friday in the Gaza Strip. This event marked peaceful protests of 40,000-50,000 Palestinian men, women and children at the border fence separating the Gaza Strip from Israel. The purpose of the Great March of Return rallies was to demanded the right of Palestinian refugees to return to cities and villages of origin from which they had been displaced by what is now called Israel (Khoury et al., 2018). They also protested in rejection to the tight blockade imposed on the Gaza Strip by the Israeli occupation. The weekly protests lasted for a year and attracting broad and diverse audiences, including men, women,

teenagers, the elderly, civil society leaders, political activists and public figures. However, the GMR activities evolved during the past twelve months to include night-time disruptions along the fence in addition to demonstrations along the Gaza coastline (MSF, 2019). At least 110 Palestinians were shot dead between 30 March to 15 May, 2018 (Sanchez, 2018).

Nowadays, sentiment analysis is gaining special emphasis due to the widespread use of social networks. Sentiment analysis is a natural language processing (NLP) field that attempt to classify and extract opinions within a document. In this respect, the phrases 'study of sentiments; and 'mining of opinions' are quite the same. Whereas the word "opinion" is broader in use than the word "sentiment," yet the two words have been used interchangeably by prior researchers. When evaluating opinions and sentiments, an audience's views and beliefs are tracked on the web to assess whether the audience perceives certain ideas or events positively or negatively. The purpose of this evaluation is intended to help companies and other concerned entities to define the various methods to enhance the quality of their goods.

With more users sharing more information across various platforms, the popularity of social media has expanded exponentially. This vast wealth of data offers unique opportunities for data mining professionals. Currently, there are over 2.50 billion active Facebook users worldwide and more than 300 million active twitter users (Clement, 2020). Fig. 1 shows the social networks users in April 2020 (in millions). Furthermore, Figure 2 shows a timeline with the number of monetizable daily active Twitter users worldwide as of the second quarter of 2020.

The political sphere and the identification of people's views and attitudes towards ongoing political activities are one of the most important areas of opinion mining. In this regard, one of the advantages of analysing political opinions is that it does not require much time and effort, and it returns acceptable results. Therefore, an experiment was designed to evaluate Arabs' opinions on the Great March of Return rallies using simple models. It is worth remembering that Arabic is one of the world's most significant languages and is used by more than 290 million people every day (UNESCO, 2012). In

this paper, four supervised learning algorithms were trained using three sentiment analysis Arabic tweet datasets. Then, the four algorithms were evaluated in terms of accuracy. Finally, the algorithm with the highest accuracy was used to predict opinion mining of Arab people on the Great March of Return rallies in the Gaza strip. The Arabic tweets were collected using social media APIs and were then preprocessed and prepared to train the most accurate algorithm.

The rest of this paper is organized as follows: Section 2 reviews some opinion mining techniques. Section 3 explains the methodology, section 4 demonstrates the experiments and results, and the conclusion is presented in Section 5.

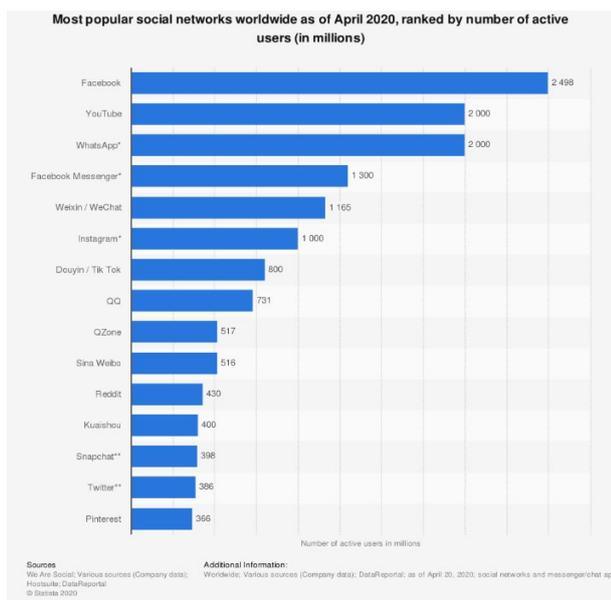


Figure 1. Social Networks as of April 2020 (in millions) (Clement, 2020)

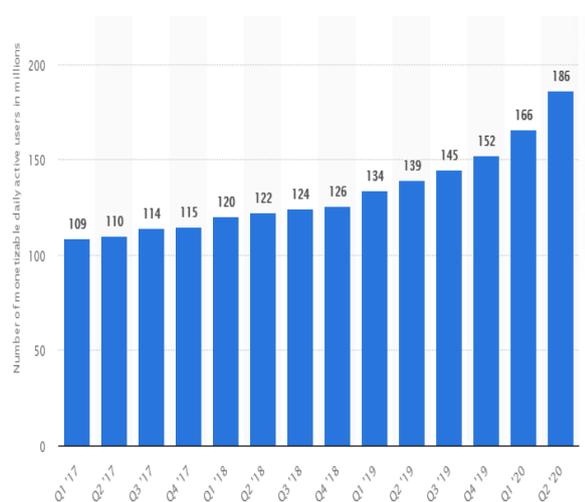


Figure 2. Number of Daily Active Twitter Users 2017-2020 (in millions) (Clement, 2020)

II RELATED WORKS

The related works are divided into three categories which are, sentiment analysis on Arabic tweets, sentiment analysis on English tweets and analyze reactions on social media about Palestinian and Israel conflict. The works achieved by classifying opinions about Arabic tweets were based on opinion mining approaches such as machine learning, semantic orientation and deep learning. A classification model was proposed by Pak & Paroubek (2010) to identify tweets as objective, positive and negative. They created an Arabic Twitter corpus by gathering tweets using the Twitter API and by automatically annotating those tweets. Using that corpus, they created a sentiment classifier that uses features like n-gram and POS-tags, based on the multinomial Naive Bayes (NB). Rushdi Saleh et al. (2011) used Support Vector Machines (SVM) method to train systems for testing different domains of data sets with numerous features. Also, the SVM-based framework was designed for subjectivity and emotion analysis for Arabic social media genres (Abdul-Mageed et al., 2012) for both Modern Standard Arabic and dialectal Arabic. The results indicated that solutions for each domain and task should be created. In addition, Omar et al. (2013) used an ensemble of machine learning classifiers; NB, SVM, and Rocchio classifiers to deal with the sentiment analysis of Arabic customer reviews. They found that NB algorithm outperformed the other algorithms on the basis of comparing the efficiency of the algorithms. The work of Duwairi et al. (2014) utilized three machine learning classifiers, namely, SVM, NB, and k-nearest classifier (KNN). The corpus encompassed data in Modern Standard Arabic. The best performance of NB was achieved when no filtering of stop words and no stemming were used. Moreover, an Arabic lexicon stored on the device was given by Badaro et al. (2015). The text is stemmed at first, and then the words are compared to their own existing ArSenL. In this respect, ArSenL is the first large-scale Standard Arabic sentiment lexicon available to the public (ArSenL). On a published collection of Arabic tweets, the method was tested, and an average accuracy of 67 percent was achieved. Recently, in order to predict the feeling of Arabic tweets, Heikal et al. (2018) implemented an ensemble model by integrating Convolutional Neural Network (CNN) and Long Short-Term Memory (LSTM) models. More recently, Al Omari et al. (2019) suggested a logistic regression method combined with the term frequency and inverse document frequency

(TF*IDF) for the Arabic evaluation arrangement in the Lebanon nation for administrative reviews. Their model was biased in predicting adverse sentiment reviews. In addition to that, we discuss the sentimental analysis related works on English language based on opinion mining approaches. As an example, Alrehili & Albalawi (2019) used a sentiment classification model to classify customer reviews using ensemble voting method which combined NB, SVM, random forest (RF), Bagging and Boosting. Moreover, there are subject related works that aims to analyze reactions on social media about a specific issue. As for Palestinians, social media is used as a tool to allow them to express their voice and opinion (Siapera, 2014). Since Palestinians people don't have their cultural and sociality equality. And the work of Siapera et al. (2015) explore the “proliferation of Palestine content in online spaces,” which they claim has kept Palestine's “memory and current issues alive” in the global arena since the Gaza War of 2014. Both studies illustrate the extensive role of digital media in the Israel-Gaza conflict, which can be incorporated into a larger theoretical context on the 'mediatisation' process. Deegan et al. (2018) discussed in his paper how to use twitter as a tool to represent the conflict between Palestine and Israel and how much the Palestinian side suffers. As far as the Palestinian-Israeli problem is concerned, very few studies have exploited the content of social media to catch trends or patterns relevant to the ongoing conflict (Siapera 2014, Siapera et al. 2015, Deegan et al. 2018). However, these studies focused on systematic statistical reviews, rather than on data mining or sentiment analysis. In our work, we use text mining techniques including some machine learning classifiers to predict the opinions of Arab people on the Great March of Return rallies in the Gaza strip during their trending on Twitter.

III METHODOLOGY

This section discusses people’s opinions on “The Great March of Return” using “مسيرة العودة الكبرى” twitter hashtag. Our framework consists of three main steps: data collection, preprocessing of the collected data and classifying the analyzed data into either positive or negative data. Fig. 3 demonstrates the basic overview of sentiment analysis framework.

A. Dataset for Training

To conduct our experiments, three sentiment analysis Arabic datasets (predefined data), as shown in Table 1 were used to train four machine learning

algorithms. The first dataset is AJGT, which is Arabic Jordanian General Tweets Corpus consisted of 1,800 tweets annotated as positive and negative (Dahou et al., 2019). The second dataset is ASTD, which is Arabic Sentiment Tweets Dataset from (Dahou et al., 2019). The third dataset is Twitter, using a tweet crawler, 2000 classified tweets were used on various topics such as politics and arts (1000 positive tweets and 1000 negative ones). These tweets contained opinions written in both Modern Standard Arabic (MSA) and the dialect of Jordan (Abdulla et al., 2013).

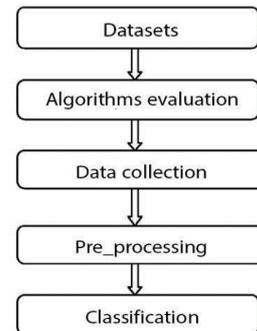


Figure 3. The Proposed Sentiment Analysis Framework

Table 1. Sentiment Analysis Arabic Datasets

Name	Positive	Negative	Total
AJGT	900	900	1800
ASTD	799	1684	2483
Twitter	1000	1000	2000

B. Machine Learning Classifiers

Machine Learning (ML) systems are categorized into three different types: supervised, unsupervised, and semi-supervised learning which combines supervised and unsupervised methods (Almunirawi & Maghari, 2016). Text classification is often performed using supervised machine learning algorithms due to the large number of labeled text datasets. The four used supervised classifiers are described as follows:

Artificial Neural Networks (ANN). ANN are a representation of a network of interconnected "neurons" that can use an objective function to calculate a set of values as inputs to generate the desired output (Irsoy & Cardie, 2014).

Decision Tree (DT). DT begins by selecting a feature as a root node, and then generates a leaf for each possible level of that feature (Trstenjak et al., 2014), (Zoroub & Maghari, 2017).

Logistic Regression (LR). LR is one of linear binary classification methods. It depends on the probability that the object belongs to a particular class (Firyulina & Kashirina, 2020). LR performs a statistical analysis to test for associations, or relationships, between variables. LR is a predictive analysis where the model

Stemming. Stemming tends to reduce a word to a popular base form in inflectional and related types, such as (أحب، يحب، يحبون، أحبوتحب).

Tokenization. Tokenization refers to the division of the sentence into its desired component elements. In all NLP activities, it is an essential step.

Feature extraction. Bag Of Word (BOW) is used for text representation. It presents the word in fixed-length, but it ignores the order of words and the grammatical structure. We also considered the TF-IDF (Term Frequency - Inverse Document Frequency), which is a popular term weighting method for feature selection (Saad & Ashour, 2010).

F. Classification

In this phase, the algorithm (classifier) with the highest accuracy score was used to predict opinion mining of Arab people on the Great March of Return rallies in the Gaza strip. The classifier was applied on the preprocessed tweets to see the opinion.

IV RESULTS & DISCUSSION

We used four different algorithms to train three different datasets and to build a model. All the training data stored in the folders (AJGT, ASTD, Twitter) were gathered in a new folder. Then, the four algorithms were simultaneously applied on the three sentiment analysis datasets. Table 3 to Table 6 show the accuracy, recall, precision, and F1-Measure values for the four classifiers.

A. Classifiers Evaluation Results

Table 3 represents the results of applying machine learning classifiers on sentiment analysis Arabic datasets. It shows that the Neural Network classifier gave an F1-Measure values of 0.82% and 74% for negative and positive tweets respectively with an average value of 78.0%. The classification accuracy was 79.0% which is considered a good value.

Table 2. Classification Report Neural Network [MLP]

	Precision	Recall	F1-Measure	Accuracy
Negative	0.77	0.89	0.82	
Positive	0.84	0.66	0.74	
Average	0.80	0.78	0.78	
Accuracy	0.79			

Table 4 represents the results of applying Decision Tree classifier on the Arabic datasets. It shows that Decision Tree classifier gave accuracy value of (73.0%) and an average F1-Measure of (71.0%).

Table 3. Classification Report Decision Tree

	Precision	Recall	F1-Measure
Negative	0.76	0.75	0.75
Positive	0.70	0.71	0.70
Average	0.71	0.71	0.71
Accuracy	0.73		

Table 5 represents the results of applying Logistic Regression classifier on the Arabic datasets. As shown, Logistic Regression Classifier achieved the best accuracy and average F1-Measure values of (83.0%) and (82.0%) respectively. Logistic Regression also returned the best average recall (82.0%).

Table 4. Classification Report Logistic Regression

	Precision	Recall	F1-Measure
Negative	0.79	0.94	0.86
Positive	0.90	0.70	0.79
Average	0.85	0.82	0.82
Accuracy	0.83		

Table 6 represents the results of applying Support Vector Machine classifier on the Arabic datasets. It shows that SVM gave an accuracy value of (82.0%) and an average F1-Measure values of 81.0%, which are considered very good results and came second in rank compared to the other three classifiers.

Table 5. Classification Report Support Vector Machine (SVM)

	Precision	Recall	F1-Measure
Negative	0.79	0.93	0.86
Positive	0.88	0.68	0.77
Average	0.84	0.80	0.81
Accuracy	0.82		

Based on the previous results, Logistic Regression achieved the best score to the predefined data. Therefore, logistic regression was applied on the collected tweets of the hashtag (#The Great March of Return) for opinion mining on GMR rallies. Dreiseitl and Ohno-Machado urged that logistic regression has lower generalization error than decision tree classifier. It also easier to build than SVM (Dreiseitl & Ohno-Machado, 2002). Bolbol & Maghari (2020) also applied some ML classifiers on Arabic tweet datasets and found that LR outperformed other classifiers as DT.

B. Opinion Mining Results on GMR

Logistic Regression classifier was applied on the 500 tweets collected from twitter using the "مسيرة العودة الكبرى" (#The Great March of Return) hashtag. A large difference was found in the results of sentiment analysis of twitter data. The classification results showed that only 71 tweets out of the 500 tweets were positive (with 14.2%), and 429 tweets were negative (with 85.8%). Overall, the number of the negative tweets was greater than that of the positive tweets.

V CONCLUSION AND FUTURE WORK

This paper aimed to introduce a simple approach for Sentiment Analysis through extracting opinions from Arabic tweets using machine learning. We used three sentiment analysis Arabic datasets to train and

evaluate four machine learning algorithms. The best accuracy achieved was 83% by using Logistic Regression. Application of logistic regression on the sample tweets revealed that 85.8% of the tweets opposed the Great March of Return, whereas 14.2% of the tweets supported it. In future, we will use CNN architecture in order to improve the results and enhance the classification performance. We can check the model over other datasets that are larger than ASTD, AJGT and Twitter, which were used primarily to equate our findings with the Arabic-language deep learning model.

REFERENCES

- Abdul-Mageed, M., Kübler, S., & Diab, M. (2012). SAMAR: a system for subjectivity and sentiment analysis of Arabic social media. 12 Proceedings of the 3rd Workshop in Computational Approaches to Subjectivity and Sentiment Analysis, July, 19–28. <http://dl.acm.org/citation.cfm?id=2392963.2392971>
- Abdulla, N., Mahyoub, N., Shehab, M., & Al-Ayyoub, M. (2013). Arabic sentiment analysis: Corpus-based and lexicon-based. Proceedings of The IEEE Conference on Applied Electrical Engineering and Computing Technologies (AEECT).
- Al Omari, M., Al-Hajj, M., Hammami, N., & Sabra, A. (2019). Sentiment classifier: Logistic regression for Arabic services' reviews in Lebanon. 2019 International Conference on Computer and Information Sciences, ICCIS 2019, 1–5. <https://doi.org/10.1109/ICCISci.2019.8716394>
- Alhaj, B. A., & Maghari, A. Y. A. (2017). Predicting user entries by using data mining algorithms. 2017 Palestinian International Conference on Information and Communication Technology (PICICT), 110–114.
- Almunirawi, K. M., & Maghari, A. Y. A. (2016). A Comparative Study on Serial Decision Tree Classification Algorithms in Text Mining. International Journal of Intelligent Computing Research, 7(4), 754–760.
- Alrehili, A., & Albalawi, K. (2019). Sentiment analysis of customer reviews using ensemble method. 2019 International Conference on Computer and Information Sciences, ICCIS 2019, 1–6. <https://doi.org/10.1109/ICCISci.2019.8716454>
- Andrew, A. M. (2001). An Introduction to Support Vector Machines and Other Kernel-based Learning Methods. In *Kybernetes* (Vol. 30, Issue 1, pp. 103–115). Cambridge university press. <https://doi.org/10.1108/k.2001.30.1.103.6>
- Badaro, G., Baly, R., Akel, R., Fayad, L., Khairallah, J., Hajj, H., Shaban, K., & El-Hajj, W. (2015). A Light Lexicon-based Mobile Application for Sentiment Mining of Arabic Tweets. Proceedings of the Second Workshop on Arabic Natural Language Processing, 18–25. <https://doi.org/10.18653/v1/w15-3203>
- Bolbol, N. K., & Maghari, A. Y. (2020). Sentiment Analysis of Arabic Tweets Using Supervised Machine Learning. 2020 International Conference on Promising Electronic Technologies (ICPET).
- Clement, J. (2020). Most popular social networks worldwide as of July 2020. <https://www.statista.com/statistics/272014/global-social-networks-ranked-by-number-of-users/>
- Dahou, A., Elaziz, M. A., Zhou, J., & Xiong, S. (2019). Arabic Sentiment Classification Using Convolutional Neural Network and Differential Evolution Algorithm. *Computational Intelligence and Neuroscience*, 2019(12), 75–85. <https://doi.org/10.1155/2019/2537689>
- Deegan, J., Hogan, J., Feeney, S., & O'Rourke, B. (2018). The Self and Other: Portraying Israeli and Palestinian Identities on Twitter. *Irish Communication Review*, 16(1), 8.
- Dreiseitl, S., & Ohno-Machado, L. (2002). Logistic regression and artificial neural network classification models: a methodology review. *Journal of Biomedical Informatics*, 35(5), 352–359. [https://doi.org/https://doi.org/10.1016/S1532-0464\(03\)00034-0](https://doi.org/https://doi.org/10.1016/S1532-0464(03)00034-0)
- Duwairi, R. M., Marji, R., Sha'ban, N., & Rushaidat, S. (2014). Sentiment analysis in arabic tweets. 2014 5th International Conference on Information and Communication Systems (ICICS), 1–6.
- Firyulina, M. A., & Kashirina, I. L. (2020). Classification of cardiac arrhythmia using machine learning techniques. *Journal of Physics*, 1479. <https://doi.org/10.1088/1742-6596/1479/1/012086>
- Frey, B. B. (2018). *The SAGE encyclopedia of educational research, measurement, and evaluation*. Sage Publications. <https://doi.org/10.4135/9781506326139.n403>
- Heikal, M., Toriki, M., & El-Makky, N. (2018). Sentiment analysis of Arabic Tweets using deep learning. *Procedia Computer Science*, 142, 114–122.
- Irsoy, O., & Cardie, C. (2014). Opinion mining with deep recurrent neural networks. EMNLP 2014 - 2014 Conference on Empirical Methods in Natural Language Processing, Proceedings of the Conference, 720–728. <https://doi.org/10.3115/v1/d14-1080>
- Khoury, Jack; Kubovich, Yaniv; Zikri, A. Ben. (2018). Mass Gaza Border Clashes. <https://www.haaretz.com/middle-east-news/palestinians/-premium-mass-gaza-border-clashes-52-killed-by-israeli-gunfire-2-410-wounded-1.6091548>
- Maghari, A. Y., & Zendah, J. H. (2019). Detecting Significant Events in Arabic Microblogs using Soft Frequent Pattern Mining. *Journal of Engineering Research and Technology*, 6(1), 11–19.
- MSF. (2019). Great March of Return. <https://www.msf.org/great-march-return-depth>
- Omar, N., Albared, M., Al-Shabi, A. Q., & Al-Moslmi, T. (2013). Ensemble of Classification Algorithms for Subjectivity and Sentiment Analysis of Arabic Customer Reviews. *International Journal of Advancements in Computing Technology(IJACT)*, 5(12), 77–85.
- Pak, A., & Paroubek, P. (2010). Twitter as a corpus for sentiment analysis and opinion mining. Proceedings of the 7th International Conference on Language Resources and Evaluation, LREC 2010, 10(2010), 1320–1326. <https://doi.org/10.17148/ijarce.2016.51274>
- Rushdi Saleh, M., Martín-Valdivia, M. T., Montejo-Ráez, A., & Ureña-López, L. A. (2011). Experiments with SVM to classify opinions in different domains. *Expert Systems with Applications*, 38(12), 14799–14804. <https://doi.org/10.1016/j.eswa.2011.05.070>
- Saad, M. K., & Ashour, W. M. (2010). Arabic text classification using decision trees. *Arabic Text Classification Using Decision Trees*, 2.
- Sanchez, R. (2018). Gaza braces for protests and funerals a day after at least 58 Palestinians killed by Israeli troops. *Telegraph*.
- Siapera, E. (2014). Tweeting #Palestine: Twitter and the mediation of Palestine. *International Journal of Cultural Studies*, 17(6), 539–555. <https://doi.org/10.1177/1367877913503865>
- Siapera, E., Hunt, G., & Lynn, T. (2015). #GazaUnderAttack: Twitter, Palestine and diffused war. *Information Communication and Society*, 18(11), 1297–1319. <https://doi.org/10.1080/1369118X.2015.1070188>
- Trstenjak, B., Mikac, S., & Donko, D. (2014). KNN with TF-IDF based framework for text categorization. *Procedia Engineering*, 69, 1356–1364.
- UNESCO. (2012). *World Arabic Language Day*. <https://en.unesco.org/commemorations/worldarabiclanguageaday>
- Zoroub, M. K., & Maghari, A. Y. (2017). Candidate Teacher Performance Prediction Using Classification Techniques: A Case Study of High Schools in Gaza-Strip. 2017 International Conference on Promising Electronic Technologies (ICPET), 129–134. <https://doi.org/10.1109/ICPET.2017.30>