

# Twitter Sentiment Analysis of Covid-19 Using Term Weighting TF-IDF And Logistic Regression

Imamah  
Faculty of Engineering  
University of Trunojoyo Madura  
Bangkalan, Indonesia  
i2m@trunojoyo.ac.id

Fika Hastarita Rachman  
Faculty of Engineering  
University of Trunojoyo Madura  
Bangkalan, Indonesia  
fika.rachman@trunojoyo.ac.id

**Abstract**—Covid-19 attack world population, and has brought much impact in all aspects of life. Stay at home and doing less in terms of social interactions. This can have a negative effect on mental health, so in this study we use sentiment analysis to know about mental health through public opinion on Twitter. Dataset which used in this study in this study is covid-19 tweets collected at 30 April 2020. Essentially, this dataset consists of 355384 tweets reviews. Covid-19 tweets will classify with the Logistic Regression method. Based on this research, the accuracy of the covid-19 tweets sentiment classification is 94.71%.

**Keywords**—coronavirus, twitter, sentiment analysis, preprocessing, positive, negative.

## I. INTRODUCTION

A sentiment analysis system usually involves of opinion mining and sentiment analysis. The prior purposes to classify the aspect and opinion revealed in reviews of data which is obtained by crawling techniques on social media or website [1],[2], and the last one to gather polarity score to classify sentiment [3],[4]. In the other words, Sentiment analysis is used to examine feelings, views, emotions, expression, beliefs, attitude, and opinion [5].

Nowadays people using social media like Twitter to express their opinion. According to *Internet Live Statistics* meanwhile 2013 the amount of tweets sent every day has stretched 500 million[6]. Some studies have been issued on the aptitude of twitter to forecast the whole thing from the stock market[7] to elections[8]. In this study, we collect data reviews from Twitter to know about public opinion in worldwide about Covid-19. As we know that Covid-19 attack World population, and has brought much impact in all aspects of life. We should stay at home and doing less in terms of social interactions. This can have a negative effect on mental health of people in the world, so in this study we use sentiment analysis to know about mental health through public opinion on Twitter.

With the sentiment analysis of covid 19 data from Twitter, we can see the mental health of people in the world, whether they still feel safe or are they in the excessive worry stage. Safe conditions can be observed from the neutral and positive sentiment values, while the worry condition can be seen from the negative sentiment values.

## II. RELATED WORK

This section will provide a literature study about sentiment analysis. Cahyo et all shown an examination of the sentiment analysis of mobile banking opinions using TF-IDF and cosine Similarity[9]. In 2019, Soumya using Malayalam

tweets dataset with machine learning techniques to analyze the sentiment. The tweets are classified into positive and negative using Support Vector Machine (SVM), Naive Bayes (NB), and Random Forest (RF) [10]. Abyan collect data reviews of the body shop tea tree oil on Twitter with positive and negative sentiments and with the accuracy level of Tea Tree Oil sentiment analysis models on female daily at 61.51%[11]. The other research performed the sentiment analysis of Go- Pay tweets. The researcher using lexicon to classify the sentiment into positive and negative label. The accuracy of linear kernel SVM is 89.17% and the polynomial kernel is 84.38% [12].

TF-IDF does not contain the number of occurrences of words alone as in BOW. But TF-IDF also looked at the important and less important words of the document. Previous research [13] compared selection features using BOW with TF-IDF. The result is that TF-IDF is better than BOW, so in this study using the TF-IDF concept in its selection features.

The choice of classification method is very important to obtain proper accuracy and in accordance with existing data. According to research [14], if the test involves positive and negative predictive values, the appropriate algorithm is binary logistic regression. In this study the predictive data are negative, positive and neutral, so we use the logistic regression method.

## III. PROPOSED METHODS

Term Weighting TF-IDF and Logistic Regression is a proposed method to classify neutral, positive and negative sentiment. The description of our proposed work as seen in the Fig. 1.

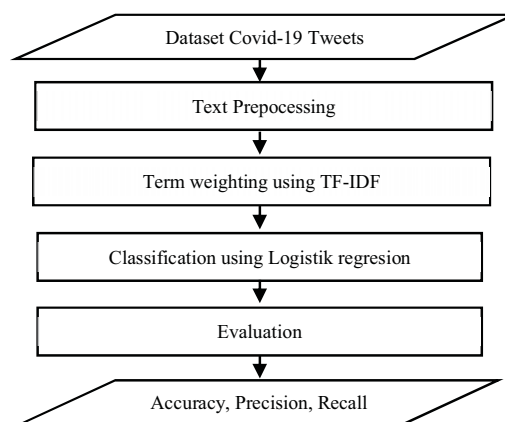


Fig. 1. Proposed methods of sentiment analysis.

In accordance with Figure 1, the initial stages are doing text preprocessing on the crawled data from Twitter. Preprocessing results in the form of terms. The term is then selected for features using TF-IDF to produce vectors. Data from the TF-IDF method were classified using the logistic regression method. Then the evaluation process is carried out with recall, precision and accuracy.

#### A. Dataset

Dataset which used in this study in this study is covid-19 tweets collected at 30 April 2020. Essentially, this dataset consists of 355384 tweets reviews. Example of dataset as seen at Fig. 2.

	created_at	text
355379	2020-04-30T23:59:58Z	dear nycmayorthe city millions spend interfait...
355380	2020-04-30T23:59:59Z	middle covid sales surge march kogan delivery ...
355381	2020-04-30T23:59:59Z	lockdown covid continues improved air quality ...
355382	2020-04-30T23:59:59Z	best answer covid wemournthemall httpstcoqrrvahgb
355383	2020-04-30T23:59:59Z	wfxtmalini cmichaalgibson boston three positiv...

Fig. 2. Covid-19 Tweets Dataset

#### B. Text Preprocessing

Text preprocessing is part of text mining to find patterns and to get valuable information for a specific purpose[15]. Based on the anomalies structure or unstructured of data text, the text processing is required some to develop data text becomes structured[16]. The steps of text preprocessing are tokenizing, filtering, and stemming.

#### C. Term Weighting TF-IDF

Term weighting is the technique that gives weights (values) to separately term in a document[17]. Three aspects influence the scheme of term weighting is term frequency(TF), inverse document frequency(IDF), and normalization.

In this study, TF- IDF is used as the feature weighting method. TF- IDF determine a weight of the term with two factors:

- Calculating the number of term  $j$  occurrences in document  $i$  or also known as term frequency denoted by  $tf_{i,j}$ .
- Calculating the frequent it arises in the entire document collection or also known as document frequency denoted by  $df_{i,j}$ .

Term weighting TF-IDF is considered formed on the following formula:

$$Tf.IDF = TF_{i,j} \times IDF_{i,j} = TF_{i,j} \times \log \frac{N}{DF_j} \quad (1)$$

where,

$N$  = number of documents in collections

TF = term frequency

IDF = inverse document frequency.

#### D. Classification using Logistic Resegion

Sentiment analysis is a process of classifying text into positive, negative and neutral. One of the method which is widely used to do supervised sentiment analysis is Logistic

regression. The classification process using this method is done by extracts real-valued features from the input, multiplies each by a weight, sums them, and passes the sum through a sigmoid function to generate a probability. A threshold value is used to make a decision.

Logistic regression able to classify sentiment into two class with positive and negative label or multiple classes using multinomial logistic regression. In this study, Logistic regression classifier is used to train and test dataset which has three class (positive, neutral and negative) as labels. A weighting scheme is adapted for separately class, which is used to adjust the probability since the training data is unbalanced. [18].

#### E. Accuracy Testing

Accuracy testing is conducted to know the performance of proposed model in classifying the sentiment into three class labels (positive, negative and netra). The measurement techniques used for evaluation the proposed method is confusion matrix, with three classes. Confusion matrix will arrange actual class results and predictive classes. The final results of this research were constructed on the precision, recall, and f-measure which are usually used in evaluation metrics in sentiment classification research. Precision measures the perfection of a classifier, and recall dealings the wholeness or compassion of a classifier. The mixture of precision and recall is measured by the f-measure, which is the weighted harmonic mean of precision and recall.

## IV. RESULT AND DISCUSSION

In this section, we the evaluation results of proposed method of sentiment classification on covid19 tweets dataset.

#### A. Preprocessing Data

Covid19 tweets dataset is unstructured and contain of noises. The dataset need to preprocessing in several stages are cleansing, eliminating numbers, emoticons, punctuation marks, case folding, filtering and tokenizing. The tools used I this research is Sklearn, Pandas and many more library of python for data preprocessing which runs on Jupyter Notebook. Example of preprocessing are shown in TABLE I.

TABLE I. PREPROCESSING DATA

Original Tweets	Preprocessing output
Microsoft sees digital reboot from pandemic, profits up #COVID19 https://t.co/ruU6qomnff	microsoft sees digital reboot from pandemic profits up covid httpstcoruqqomnff

#### B. Term Weighting

The words which is a result of preprocessing data will be converted into numeric using term weighting. This process aims to calculate the weight of each word that will be used as a feature, the more documents will be processed then the more features. At this stage there are two part of process namely TF (Term Frequency) and IDF (Inverse Document Frequency).

TF is the number of words that appear in each document. IDF is the number of document values in each word that are inversely proportional, it means if a word rarely appears in a document, the IDF value is greater than words that often appear. The result of weighting the word with TF-IDF is the multiplication of the values of TF and IDF which will produce less weight if the word often appears in every document in the

collection, on the contrary the weight of TF-IDF will be greater if the word rarely appears in each document in collection. By using TF IDF, you can also analyze the most tweeted words as in Fig. 3 and Fig. 4.

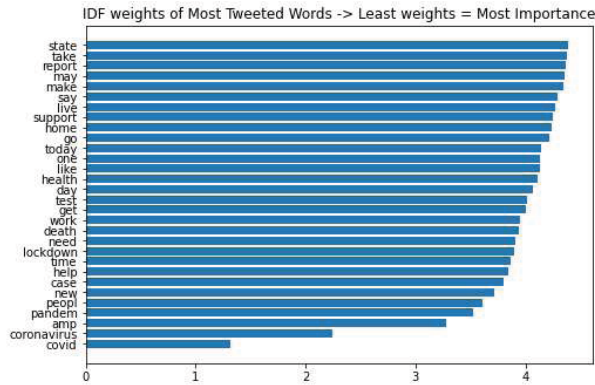


Fig. 3. IDF weights of most tweeted words

It can be seen in the graph in Fig. 3 that the terms that dominate the document in the top 10 terms are terms with neutral and positive sentiment values, for terms that have negative sentiment are below. But on the TF graph (Figure 4), it can be seen that the terms that have a "negative" sentiment like "death" are higher in value than those with a positive sentiment like "health".

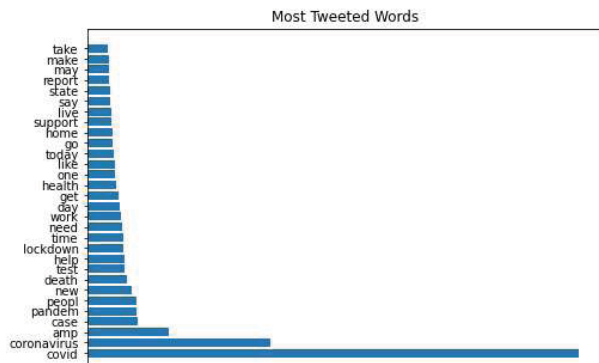


Fig. 4. Most Tweeted words using Term Frequency

### C. Data Labelling

The next steps is data labelling which purposes to give positive, neutral and negative labels on Covid-19 Tweets dataset. In this study, we used polarity score to compute score of sentiment. Covid-19 reviews will be positive labelled if a score larger to 0. Otherwise if a score is less than to 0, it will be negative labelled. Covid-19 reviews will be neutral labelled when the score is equal to 0. The source code of sentiment polarity is shown in Fig. 5.

```
polarity_score = []

for text in df['text']:
    blob = TextBlob(text)
    for sentence in blob.sentences:
        if sentence.sentiment.polarity > 0.0:
            polarity_score.append('positive')
        elif sentence.sentiment.polarity < 0.0:
            polarity_score.append('negative')
        else:
            polarity_score.append('neutral')
```

Fig. 5. Source code to find polarity score fo sentiment

Total covid-19 tweets dataset is 355384 which includes 213122 neutral sentiment, 100683 positive sentiments and 41579 negative sentiments. Fig. 6 summarizes the statistics of this dataset.

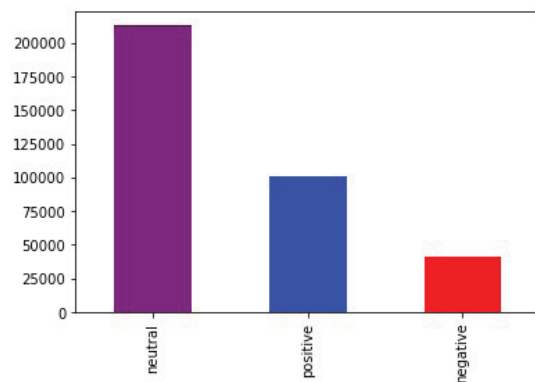


Fig. 6. Sentiment Classification of Covid-19 Tweets Dataset.

### D. Data Visualization using Wordcloud

Wordcloud is used to display the words that often appear in the document. Here is the result of wordcloud visualization for covid-19 tweets. Fig. 7. describes the words that most appear which has positive label. The Examples of the most common words health, good, need, great, help.



Fig. 7. Wordcloud of positive sentiment.

Meanwhile Fig. 8 shows the words that most appear which has negative label, for example are cases, lockdown, bad, pandemic.



- [11] K. D. Adhine Salsabila, A. Ghiffarie, R. P. Baistama, M. I. Variadi, and M. D. Rhajendra, "Sentiment Analysis of The Body Shop Tea Tree Oil Products (*Analisis Sentimen Terhadap Produk The Body Shop Tea Tree Oil*)," *J. Teknol. dan Manaj. Inform.*, vol. 5, no. 2, 2019.
- [12] R. Mahendrajaya, G. A. Buntoro, and M. B. Setyawan, "Gopay User Sentiment Analysis Using the Lexicon Based Method and Support Vector Machine (*Analisis Sentimen Pengguna Gopay Menggunakan Metode Lexicon Based dan Support Vector Machine*)," *Komputek*, vol. 3, no. 2, p. 52, 2019.
- [13] V. L. Nguyen, D. Kim, V. P. Ho, and Y. Lim, "A New Recognition Method for Visualizing Music Emotion," in *International Journal of Electrical and Computer Engineering (IJECE)*, 2017, vol. 7, no. 3, pp. 1246–1254.
- [14] F. Zabli and I. H. Osman, "ReviewModus : Text classification and sentiment prediction of unstructured reviews using a hybrid combination of machine learning and evaluation models," *Applied Mathematical Modelling*, vol. 71, no. 2019, pp. 569–583, 2020.
- [15] I. Imamah, H. Husni, E. M. Rohman, I. O. Suzanti, and F. A. Mufarroha, "Text mining and Support Vector Machine for Sentiment Analysis of tourist Reviews in Bangkalan Regency," *J. Phys. Conf. Ser.*, vol. 1477, no. 2, pp. 0–6, 2020.
- [14] O. Somantri and D. Dairoh, "Sentiment Analysis of Tegal City Tourism Destination Assessment Based on Text Mining (*Analisis Sentimen Penilaian Tempat Tujuan Wisata Kota Tegal Berbasis Text Mining*)," *J. Edukasi dan Penelit. Inform.*, vol. 5, no. 2, p. 191, 2019.
- [15] H. Wu and X. Gu, "Reducing over-weighting in supervised term weighting for sentiment analysis," *COLING 2014 - 25th Int. Conf. Comput. Linguist. Proc. COLING 2014 Tech. Pap.*, pp. 1322–1330, 2014.
- [16] H. Hamdan, P. Bellot, and F. Bechet, "Lsislif: CRF and Logistic Regression for Opinion Target Extraction and Sentiment Polarity Analysis," no. *SemEval*, pp. 753–758, 2015.