# AUTOMATIC GENERATION OF WORD-EMOTION LEXICON FOR MULTIPLE SENTIMENT POLARITIES ON SOCIAL MEDIA TEXTS

LIZA WIKARSA[1] AND MINSOO KIM[2,*]

[1]Graduate School of Management of Technology
Pukyong National University
Yongdang Campus, 365 Sinseon-ro, Nam-gu, Busan 48547, Korea
liza_wikarsa@pukyong.ac.kr

[2]Division of Systems Management and Engineering
Pukyong National University
Daeyeon Campus, 45 Yongso-ro, Nam-gu, Busan 48513, Korea
*Corresponding author: minsky@pknu.ac.kr

ABSTRACT. *The applications of sentiment analysis (SA) on social media texts have attracted much attention in both public and commercial sectors. Due to the limitations of general lexicons as resources for evaluating the sentiment of a text passage, this study proposes a framework to generate an automated word-emotion lexicon on Twitter by searching relations between words and tweets. The dictionary- and corpus-based approaches are utilized to classify words and tweets into seven sentiment polarities while applying the linguistic features like "and, but, or". Normalized distribution is used to calculate sentiment score for each word or tweet. The effectiveness of the sentiment classification related to health, technology, and education topic is tested. The results indicate that accuracy and precision rates are higher at the word-level than at the tweet-level by identifying a linear association between candidate words and seed words. This expanded lexicon is useful for performing domain specific tasks on social media texts.*
**Keywords:** Lexicon, Sentiment analysis, Emotion, Data mining, Social media texts

1. **Introduction.** Emotion affects individual behaviors which contribute to one's decision makings and interaction with other people. Nowadays, people eagerly express their emotions on various social media sites that result in an increase in the amount of user-generated content over time. To date, Twitter has 330 million of monthly active users while the total number of tweets sent per day is 500 million [1]. Hence, social media is seriously considered as a gold mining for human opinions which consequently raises the interest in doing in-depth research and business around the sentiment analysis (SA) field.

SA refers to an automatic extraction of opinions or sentiments that provide both public and commercial sectors with efficient tools and solutions for tracking the public opinions on their product and service, political strategy planning, financial forecasting, and more [2,3]. Marquez et al. argued that many SA methods heavily rely on opinion lexicons as resources for evaluating the sentiment of a text [4]. A sentiment lexicon is a collection of words or phrases that are commonly used to express emotion. General Inquirer (GI), SentiWordNet (SWN), and Opinion Lexicon (OL) are several examples of general lexicons (GLs) [4-6]. GLs mostly classify sentiment words into positive, negative, and neutral only. Unfortunately, GLs are unable to capture the diversity and sparseness of informal expressions used in social media texts [4-9]. Should one use a GL for domain-specific tasks, severe misunderstanding of the information can happen due to its inability to properly

match the domain in which the lexicon that was built or try to do particular tasks. Many GLs are lack of words related to specific topics. Also, the manual creation of a sentiment lexicon is really time consuming and costly [4]. In terms of the sentiment scoring, imposing more weights on certain words and assigning higher weights to the scores obtained from matching adjectives and adverbs are few methods used in the previous studies which are not fully explained and considered as subjective scoring [5].

This study proposes a framework to generate an automated word-emotion lexicon on social media texts using Eikman's six sentiment classification that are happiness, sadness, anger, surprise, fear, and disgust [2,8,10]. In addition, 'neutral' is added for non-sentiment words. Data is collected from Twitter and AFINN is opted for the seed lexicon. This framework uses a hybrid approach to integrate the dictionary- and corpus-based approaches that incorporate linguistic features to train a corpus of emotion-annotated tweets in the seed lexicon and use the model on unlabeled tweets with a probability normal distribution of the seven sentiment polarities. Thus, the framework will automatically build or expand the seed lexicon. Three SA tasks performed for evaluation are: 1) at the word level, it identifies the sentiment polarity and score for individual words based on the rules, 2) at the tweet-level, it identifies the overall sentiment of individual tweets, and 3) at the topic level, it assesses sentiments of Twitter's users with regard to health, technology, and education. The main contributions of this research are: 1) automatically generate and expand the seed lexicon that contains seven sentiment polarities which are more informative and adaptive; 2) sentiment scoring to classify user opinions objectively; 3) the proposed framework is suitable for supporting implementation of domain specific lexicon-based applications, especially on social media texts.

The rest of this paper is organized as follows. Section 2 outlines the previous studies related to the SA and lexicon. The proposed method is described in Section 3; meanwhile, the experimental results and discussion are presented in Section 4. Section 5 provides conclusions and recommendations for the future work.

## 2. Related Works.

2.1. **Sentiment analysis.** The automatic extraction of opinions for SA can be done using machine learning (ML) and lexicon-based approach (LBA) [5,6]. According to Deng et al., supervised cases normally follow the ML approach where the sentiment detection task is considered as a classification problem. ML sentiment classifiers perform well on the domain/genre of training but poorly on a different domain/genre like the social media that are diverse and constantly change [6]. The power of ML resides in the training data but it requires efforts and time to generate high-quality training data [5]. This poses challenges for analyzing social media texts due to the constant evolution of the language and increased online in user-generated content. Meanwhile, an LBA is often used in unsupervised cases where one can adopt dictionary-based, corpus-based, or both approaches [5] which will be discussed in the next section. The differences between ML and LBA in terms of their basis, algorithms, limitations, and applications can be referred to [1,2,5-7].

2.2. **Lexicon-based approach.** An LBA basically relies on a pre-defined sentiment lexicon to determine the general sentiment polarity of a given document by exploiting the word presence in the document and/or leveraging the existing lexicons containing polarized or emotional words [1,2]. Though existing lexicons are useful for many tasks, they are fixed resources that need to be improved when the domain changes. Hence, it is important to construct new lexicons for: 1) capturing different dimensions of the sentiment for specific tasks [1], 2) being sensitive to the norms of specific domains [7], 3) much larger lexicons which can be developed inferentially [8].

There are two automatic LBAs such as the dictionary- and corpus-based approaches [5,11,12]. The former is used when new sentiment words (henceforth called CW for candidate word) are identified through their relationship with the seed lexicon by consulting the dictionary for their synonyms and antonyms [5]. Liu [8] outlined the following steps for this approach: 1) a small set of sentiment words with known positive or negative orientations is first collected manually, 2) the algorithm grows this set by searching in the WordNet or other online dictionaries for synonyms and antonyms, 3) the newly found words are added to the seed list and the next iteration begins, and 4) the iterative process ends when no new words are found. To determine the sentiment polarity, several bootstrapping methods were proposed including Markov random walk model and pointwise-mutual information (PMI). However, the dictionary based approach is unable to find opinion words with domain and context specific orientations [9]. The corpus-based approach is used when new CWs are recognized based on their mutual relationships [2]. It exploits co-occurrence patterns of words found in the unstructured textual documents [3]. Liu [8] demonstrated two scenarios of this approach: 1) given general-purpose sentiment words, discover other sentiment words and their polarities from a domain corpus, and 2) adapt a GL to a new one using a domain corpus for SA applications. A set of linguistic rules can be used to identify sentiment words and their orientations from the corpus, like the use of 'and' for conjoined adjectives that usually have the same orientation. For instance, "this house is beautiful and spacious", the 'beautiful' has the positive sentiment, it can be inferred that 'spacious' is also positive. Other connectives like 'or', 'but', 'either-or', 'neither-nor' can be applied for two conjoined adjectives that have opposite sentiment polarities. This study adopts these two rules by applying: 1) conjunction 'and' for assigning the same sentiment polarity, and 2) 'but' and 'or' are for the opposite connectives.

The limitations of LBAs in the previous studies include subjective sentiment scoring of opinion words, limited coverage of domain specific words, and only three sentiment classes are classified. The effectiveness of the existing lexicons also poses challenges due to the evolving textual social media data that can result in low accuracy of the classifier in the SA of online contents. To address these issues, this study will implement a rule-based scheme using an improved version of dictionary- and lexicon-based approaches by which domain specific vocabulary is introduced to improve the efficacy of classifying the seven distinct sentiment polarities. A probability normal distribution is also used for sentiment scoring of opinion words. This study attempts to reduce data sparseness and incorrect classification of opinion words related to health, technology, and education topics.

3. **Proposed Method.**

3.1. **Data collection.** Information sources used in this framework include sentiment-annotated tweets, unlabeled tweets, and a hand-annotated seed lexicon modified from AFINN. AFINN includes slang, obscene words, acronyms and Web jargon [3]. These sources enable the model not only to incorporate precedent knowledge from the existing seed lexicon and sentiment-annotated tweets but also learn the sentiment of words in tweets that are not necessarily classified according to the sentiment polarities used in this study. Table 1 enlists 2,511 words spread through the seven polarities in the seed list. For training data, each sentiment polarity has 1,000 tweets. As for testing, this study assesses the opinions of Twitter's users on health, technology, and education. Health and technology are two out of ten topics that are mostly engaged by Twitter's users [10]. Whilst, education has significant effects on how one thinks, feels, and acts in a daily life. Each topic will have 5,000 tweets to test at the word-, tweet-, and topic-levels.

3.2. **Sentiment scoring.** A sentiment value for a sentiment polarity is scored on a seven-point scale (1 = the lowest and 7 = the highest). Then, this value is assigned to individual

TABLE 1. Number of words in the seed lexicon

| Sentiment polarity | Happiness | Surprise | Neutral | Sadness | Fear | Anger | Disgust |
|---|---|---|---|---|---|---|---|
| Number of words | 537 | 76 | 374 | 299 | 611 | 303 | 311 |

TABLE 2. Sentiment scoring

| | Disgust | Anger | Fear | Sadness | Neutral | Surprise | Happiness |
|---|---|---|---|---|---|---|---|
| Strength | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
| Normalized $f(x_i)$ | 0.082 | 0.177 | 0.322 | 0.500 | 0.678 | 0.823 | 0.918 |
| **Scoring** | **0.024** | **0.051** | **0.092** | **0.143** | **0.194** | **0.235** | **0.262** |

polarity using the normalized distribution as is given in Equation (1). Afterwards, the value of each normalized sentiment polarity is divided by the sum of the normalized values of the seven sentiment polarities as presented in Table 2. The sum of all the sentiment scores is equal to one.

$$f(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(\mu-x)^2}{2\sigma^2}}, \ \mu = \text{mean}; \ \sigma = \text{standard deviation} \tag{1}$$

$$\text{Sentiment score} = f(x_i) \left/ \sum_{i}^{n} f(x_i) \right. \tag{2}$$

3.3. **Method description.** This framework consists of three main stages, namely pre-processing, processing, and evaluation as shown in Figure 1. The procedure is done on



sw = a seed word in the seed lexicon

* Words 'and, but, or' are not removed from the stopwords list. These words are used in the rules to determine two conjoined words whether they have the same or opposite polarity.

** There are three possible conditions checked before doing the word pairing and calculating the occurrence of the pairing words.

FIGURE 1. The proposed framework

word-, tweet-, and topic-levels. For the processing stage, it involves the word- and tweet-levels. Whilst, the evaluation stage will be performed on the three levels.

**Stage 1: Pre-processing.** The total of 7 steps involved in this stage are explained in Figure 1. At the end of this stage, a list of CWs will be generated as inputs for the next stage.

**Stage 2: Processing.** At the word level, it will find the relationships between the new sentiment words (i.e., CW) and the seed words belong to each polarity in the seed list. Three conditions are checked in relation to the presence of the CWs: 1) a single CW, 2) two CWs joined using 'and', and 3) two CWs joined using 'or' or 'but'. For the second and third conditions, one of the CWs must have already existed in the seed list. The iterative processes end when no more candidates in the candidate list. Next, it compares the frequency of occurrence for each sentiment polarity. At the tweet-level, it searches for sentiment words and counts their frequency of occurrence in individual tweet. The most frequent occurrence of the word will be the dominant sentiment polarity for that tweet. At the topic level, the overall score is calculated for all the seven sentiments for each topic. The highest score of the sentiments will be the dominant sentiment for that topic. The output of this stage is an expanded lexicon with additional new words along with its sentiment polarity and score.

---

1) For a single CW:
   a. Use the dictionary-based approach to find the relationship between the CW and the seed words in each sentiment polarity.
   b. If the CW already exists in the seed list, this CW will be removed from the candidate list.
   c. If not exist yet, it will pair the CW with each seed word for every sentiment polarity. Next, it counts occurrence of the pairing word in the seed list. Once all the words in that sentiment polarity have been paired up with the CW, it will sum the frequency of occurrence for that particular sentiment polarity. The sentiment with the highest frequency of occurrence will be assigned to the CW and the sentiment score of the candidate is referred to Table 2.

*$2^{nd}$ and $3^{rd}$ conditions will use the linguistic rules from the corpus-based approach.*

2) For two CWs joined using 'and':
   If one of the words has already existed in the seed list, it will therefore automatically assign the sentiment polarity and score of the seed word to the CW. For example, 'amazed' **and** 'pleased' words whereby 'pleased' belongs to the happiness polarity in the seed lexicon. Since the 'and' conjunction word is used, the word 'amazed' (the new candidate) will automatically be assigned to the same polarity and score as the word 'pleased'.

3) For two CWs joined using 'or' or 'but':
   When two words are joined using 'but' or 'or', the sentiment polarity and score of the candidate will take up the opposite sentiment polarity of the existing word. For instance, consider 'nice' but 'expensive' expression, the word 'nice' is already in the seed lexicon and belongs to the happiness polarity. Therefore, the word 'expensive' will have the 'Sadness' polarity instead.

---

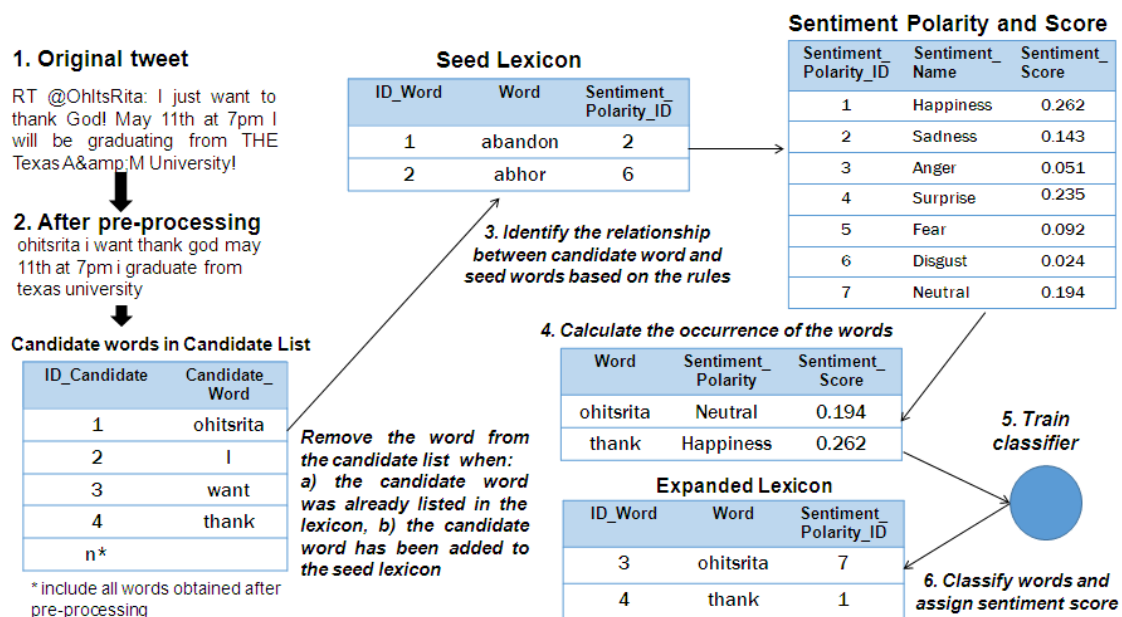FIGURE 2. Three conditional processing at the Stage 2

**Stage 3: Evaluation.** In the last stage of the evaluation, the model will be run on the test datasets to assess its performance under three different datasets, namely health, technology, and education, from Twitter. Accuracy, precision, recall, and F-measure are used to evaluate the performance of the lexicon applied on the datasets. Evaluation measures are estimated using 5 times 10-fold-cross-validation.

4. **Results and Discussion.** In the processing stage, the word-level classifier is trained using the word-level features and the word labels from the seed list. Thus, the trained classifier identifies and classifies the polarities and scores for the CWs in the candidate list, either by removing the CWs when the conditions are not being satisfied or adding to the CWs to the seed list. At the tweet level, some tweets are manually labelled and

they are used as the baseline for training. In this level, it searches for sentiment words and counts their frequency of occurrence by which the most frequent occurrence of the word will be the dominant sentiment polarity for that particular tweet. For example, there are several tweets containing a variety of sentiment words and occurrences:

1. Happiness (3), surprise (2), and neutral (1) $\longrightarrow$ the sentiment polarity is happiness because the most frequent sentiment here is happiness.
2. Fear (2), disgust (2), and neutral (1) $\longrightarrow$ the sentiment polarity is fear because fear has a higher sentiment score despite the fact that fear and disgust have the same number of occurrences.

Lastly, when all the tweets have been assigned with appropriate sentiment polarities, it will thus compare the frequencies among polarities for each topic. The sentiment polarity with the highest frequency will be the ultimate sentiment for that particular topic along with its sentiment score. The implementation of the proposed framework is shown in Figure 3.



To calculate the performance of the classifier, these sentiment polarities are gathered as follows.

1. Happiness, surprise, and neutral belong to the positive sentiment classes.
2. Sadness, anger, fear, and disgust are negative sentiment classes.

FIGURE 3. Implementation of the proposed framework

In the evaluation stage, there are 35,000 words spread across seven polarities with 15,000 tweets on the testing datasets. Neutral words are considered as positive sentiments to increase the probability of classification in this study. After processing the datasets on both word- and tweet-levels, it is interesting to find out that neutral words are the most frequent words in both datasets. However, the number of neutral words will decrease as the training data increases. Precision, recall, and F-measure are used to evaluate the performance of the proposed method in classifying these tweets on the dataset at the word-, tweet-, and topic-levels as depicted in Table 3 and Table 4.

At the word-level, the classifier achieves better results due to the use of accurately labelled seed words and a linear association between words and the sentiment of the tweets in which they occur. Experiments were conducted to classify the same words using the seven polarities and conventional methods of two classes. The results show that the positive class words going to the happiness, surprise, and neutral polarities are 53.53%, 12.93%, and 33.53% respectively. Whilst, the negative class words belonging to the fear, disgust, angry, and sadness polarities were 38.54%, 22.01%, 19.84%, and

TABLE 3. Testing at word- and tweet-levels

|  | Word-level | Tweet-level |
|---|---|---|
| **Precision** | 0.85 | 0.79 |
| **Recall** | 0.96 | 0.90 |
| **F-measure** | 0.91 | 0.87 |
| **Accuracy** | 0.84 | 0.80 |

TABLE 4. Topics and sentiment polarities

| **Topic** | **Sentiment polarity** | **Score** |
|---|---|---|
| **Health** | Neutral | 0.91 |
| **Technology** | Neutral | 0.93 |
| **Education** | Neutral | 0.95 |

19.6% respectively. If the proposed approach is to be used on any other GLs that only employ two or three sentiment classes, then constructing sentiment rules and an in-depth training are required in order to learn the relationship between these sentiment classes and our seven polarities. In addition, GLs and the proposed approach use different scoring techniques which consequently require a computational conversion to match the scores. At the tweet-level, the result is not as good as expected since it is more complex to associate and contextualize all words in a tweet to determine the dominant sentiment for the tweet. At the topic-level, the findings indicate that Twitter's users seemed to have mostly neutral opinions on these three selected topics which were a surprise as they normally express strong opinions on other topics.

Saif et al. [7] pointed out that LBA does not require training data as one can use GLs along with their sentiment classes. In contrast, this study employs seven polarities with different sentiment scoring that does not exist yet. So, this particular lexicon requires in-depth model training to identify and classify the words into the correct polarities. Contrary to Deng et al.'s belief that the LBA is suitable for real-time sentiment classification given its relatively lower computation requirement [5], this study found that the training and testing processes are time consuming and thus need higher computation requirements. As there are more CWs in the candidate list, pairing searches take longer time with the sentiment rules elaborated in Section 3. Deng et al.'s research only used three sentiment classes and fewer tweets to train and test that compelled less computation loads and processing time. To reduce these requirements, the classifier needs to recognize the context and meaning of the words. The $n$-gram model is useful for minimizing the processing time, where $n$ refers to the number of grouped words. Creating a vocabulary of word pairs is also preferable for classification tasks, rather than in isolated individual words, stressing on the importance of learning the patterns of collocations to avoid errors and increase the proficiency in using the model on unlabeled tweets. Furthermore, LDA bag of words (BoW) model can be utilized not only to uncover "the ontological structuring of the knowledge inside the individual retrieved documents in terms of words, concepts and topics" [16] but also to get the highest probability for the dominant sentiment of individual tweets. This computational overload issue needs further investigation.

The expanded lexicon can provide better improvements in performance over the seed lexicon as there are more words collected along with their sentiment polarities and scores. The results showed that the proposed framework could achieve higher precision and recall rates when neutral words are considered and treated as positive sentiments. This framework is useful for expanding the lexicon by classifying words and entire tweets to appropriate sentiment polarities and classes, showing more improvement in performance than the original seed lexicon.

5. **Conclusions.** This study is able to generate a domain-specific sentiment lexicon using the hybrid approach in the proposed framework. This lexicon can be used and applied in many domains and is able to overcome the diversity and sparseness of informal expressions used in social media texts. The effectiveness of the generated sentiment lexicon can be significantly improved that one should continuously train it and expand the seed lexicon by adding more new words and their sentiment polarities.

For the future work, it is suggested to include POS-tagged, lemmatization, and negation features to benefit from the knowledge based on their word definition and contexts. It can further increase the chance of adding non-sentiment words to the existing lexicon. Another recommendation is to not remove punctuation as it provides more clarity and stress in sentences. There is a need to add more linguistic rules, such as "either-or", "neither-nor", to better identify and classify the sentiment polarity. Lastly, it is recommended to use word collocations and LDA BoW model in order to perform better training and testing on the model. It can uncover the ontological structuring of the knowledge inside individual tweet in terms of words and topics.

**REFERENCES**

[1] E. S. Tellez, S. Miranda-Jimenez, M. Graff, D. Moctezuma, O. S. Siordia and E. A. Villaseñor, A case study of Spanish text transformations for Twitter sentiment, *Expert Systems with Applications*, vol.81, pp.457-471, 2017.

[2] M. Giatsoglou, M. G. Giatsoglou, K. Diamantaras, A. Vakali, G. Sarigiannidis and K. C. Chatzisavvas, Sentiment analysis leveraging emotions and word embeddings, *Expert Systems with Applications*, vol.69, pp.214-224, 2017.

[3] Q. Wang, Y. Jin, T. Yang and S. Cheng, An emotion-based independent cascade model for sentiment spreading, *Knowledge-Based Systems*, vol.116, pp.86-93, 2017.

[4] F. B. Marquez, E. Frank and B. Pfahringer, Building a Twitter opinion lexicon from automatically-annotated tweets, *Knowledge-Based Systems*, vol.108, pp.65-78, 2016.

[5] S. Deng, A. P. Sinha and H. Zhao, Adapting sentiment lexicons to domain-specific social media texts, *Decision Support Systems*, vol.94, pp.65-76, 2017.

[6] A. Muhammad, N. Wiratunga and R. Lothian, Contextual sentiment analysis for social media genres, *Knowledge-Based Systems*, vol.108, pp.92-101, 2016.

[7] H. Saif, Y. He, M. Fernandez and H. Alani, Contextual semantics for sentiment analysis of Twitter, *Information Processing and Management*, vol.52, pp.5-19, 2016.

[8] B. Liu, *Sentiment Analysis and Opinion Mining*, Toronto, Morgan & Claypool, 2012.

[9] W. Medhat, A. Hassan and H. Korashy, Sentiment analysis algorithms and applications: A survey, *Ain Shams Engineering Journal*, vol.5, no.4, pp.1093-1113, 2014.

[10] F. Colace, L. Casaburi, M. De Santo and L. Greco, Sentiment detection in social networks and in collaborative learning environments, *Computers in Human Behavior*, vol.51, pp.1061-1067, 2015.

[11] H. G. Liu and J. H. Guan, A model of fuzzy normal distribution, *Open Journal Statistics*, vol.6, pp.749-755, 2016.

[12] L. Wikarsa and S. N. Thahir, A text mining application of emotion classifications of Twitter's users using Naïve Bayes method, *Proc. of the 1st International Conference on Wireless and Telematics (ICWT 2015)*, Manado, Indonesia, 2015.

[13] L. Williams, C. Bannister, M. Arribas-Ayllon, A. Preece and I. Spasic, The role of idioms in sentiment analysis, *Expert Systems with Applications*, vol.42, pp.7375-7385, 2015.

[14] A. Nanji, *The Most Popular Topics on Facebook and Twitter*, http://www.marketingprofs.com/charts/2014/25346/the-most-popular-topics-on-facebook-and-twitter, accessed 22nd May 2018.

[15] F. A. Pozzi, E. Fersini, E. Messina and B. Liu, *Sentiment Analysis in Social Networks*, Cambridge, Elsevier Inc., 2017.

[16] P. D. Rocca, S. Senatore and V. Loia, A semantic-grained perspective of latent knowledge modeling, *Information Fusion*, vol.36, pp.52-67, 2017.