# Comparison of two methods on vector space model for trust in social commerce

**Hla Sann Sint, Khine Khine Oo**
University of Computer Studies, Yangon, Myanmar

| Article Info | ABSTRACT |
|---|---|
| | The study of dealing with searching information in documents within web pages is information retrieval (IR). The user needs to describe information with comments or reviews that consists of a number of words. Discovering weight of an inquiry term is helpful to decide the significance of a question. Estimation of term significance is a basic piece of most information retrieval approaches and it is commonly chosen through term frequency-inverse document frequency (TF-IDF). Also, improved TF-IDF method used to retrieve information in web documents. This paper presents comparison of TF-IDF method and improved TF-IDF method for information retrieval. Cosine similarity method calculated on both methods. The results of cosine similarity method on both methods compared on the desired threshold value. The relevance documents of TF-IDF method are more extracted than improved TF-IDF method.<br><br>*This is an open access article under the CC BY-SA license.* |

*Corresponding Author:*

Hla Sann Sint
University of Computer Studies
Yangon, Myanmar
Email: hlasannsint@ucsy.edu.mm, hlasannsint@gmail.com

## 1. INTRODUCTION

Information retrieval (IR) is answerable for capacity and retrieval of large amounts of data in a productive way. Most IR structures figure a numeric score on how well each article in the database orchestrate the ask, and rank the things as showed by this worth. The top arranging articles are then appeared to the user. The strategy may then be iterated if the user wishes to refine the request. Objective of IR is to find records suitable to data need from a huge review set. A significant deficiency of retrieval pattern information recuperation procedure is that the language that searchers use is as often as possible not equal to the one by which the information has recorded. An enormous part of the current artistic information recovery approaches depends upon a lexical match between words in customer's sales and words in target objects.

In this methods, archives and inquiries has spoken to as vectors, with every part inside the vector taking on a worth dependent on the nearness or non-attendance of a word inside the content. To decide the importance of a record for a given question, a similitude activity (ordinarily a spot item) led on the vectors yielding a solitary number. In vector space model, feature is weighted by using numbers, some commonly used with weighting methods, such as weighed of boolean, weighted of frequency, weight of term frequency-inverse document frequency (TF-IDF), weight of term frequency in class (TFC), weight of log-weighted term frequency (LTC), entropy weighting. TF-IDF weighting is the most commonly used one among them [1, 2].

## 2.　RELATED WORKS

The focus of must be researched in retrieval of information is the retrieval of data from unstructured data sources. Intuitively, this is a much harder problem than structured data retrieval. Usually, the Boolean retrieval model relegates 1 or 0 subject to the proximity or non-appearance of the terms in a document. This model performs tragically in addressing for a document. A short time later, the vector space model introduced for situated retrieval. It is broadly utilized in questioning documents, bunching, arrangement and other data recovery tasks since it is basic and straightforward.

Nowadays extraordinary trust figuring methods are available. Many complexities enrolling methods for trust have been proposed by researchers from different ways; by the by, most by a wide margin of them essentially measure certain trust related course and go along with them into a trust a spark by translating a store for each course. Different works exit on the issue of trust estimation (some may decipher it as trust esteems or trust computing) [2-4].

We-Intention, moral trust and self-motivation [5] depends on inspirational components. It has been checked by Cronbach's α, squared multiple correlations (SMC). Cronbach's α ascertain the estimation of credited to confide in score change that is an extent of the difference in the test score. The central work is to take a gander at and separate the rule segments of influencing data sharing development in social joint exertion.

Lu S. [1] proposed an improved approach method tf-idf IG to remedy this defect by information gain from information theory. This paper conquers is the restriction of old tf-idf. The idf can't well show discriminative and significance of highlight, weight change strategy is advanced in which the IDF work supplanted by assessment function used in feature selection.

A social context-aware trust sub-network extraction model processes to discover close ideal arrangements viably and proficiently utilizing ant colony algorithm (ACA) heuristic techniques [6]. The gold of this outcome extricated sub-systems inside the comparable execution time. An improved direct trust assessment technique dependent on the leader-follower clustering algorithm for context-aware trust model [7]. It attempted to tackle the issue of information sparsity issue brought about by the decent variety of services and contexts.

A collaborative filtering algorithm calculates dependent on network factorization and multi-way trust degree combination. It utilized a blend of strategy the lattice decay procedure and informal communities trust model [8]. The fundamental work is the issue of low proposal exactness brought about by the high conditionality of information. Different strategies predict the trust score of source user on track user by proliferating. Expected to discover various calculations to compute the most confided in way in companion of companion additionally in as least time-weight and credit assessment [9].

## 3.　BACKGROUND THEORY

The two essential IR models utilized during the time spent retrieving data: Boolean, and probabilistic. Likewise, these two ways attempt to order archives dependent on their pertinence to the client's needs. Among the most mainstream of approaches, because of straightforwardness and relative viability, are vector put together these ways based with respect to a plan using the recurrence of words in the record and report assortment. One of these plans is known as tf-idf. While famous, tf-idf plans may additionally improve the viability of this method [10-12].

### 3.1.　Term frequency-inverse document frequency (TF-IDF) method

TF-IDF is a customary methodology which is utilized to discover the term significance by discovering weight of a term. Steps to discover weight of a question and terms in web reports utilizing vector space model are as per the following:
−　Removed tag for desire page.
−　Remove stopwords.
−　Tokenized the given sentences.
−　Apply with poster stemming calculation.
−　Term frequency (TF) calculation for each review (query) (q) within a query (Q) from (web page) document.
−　Inverse document frequency (IDF) calculation of each term in the web pages (query (Q))
−　Compute TF-IDF of each term of query using (1) and (2).

Term recurrence (TF) is basically a rate signifying the number of times a word appears in a document. It is numerically communicated as appeared in (1). Where, $n_i$ is the number of occurrences of the considered term and $n_{max}$ is the count of the term with maximum occurrences in the web page. Inverse document frequency (IDF) considers that numerous words happen commonly in numerous archives. IDF is numerically communicated as appeared in (2). Where, $D$ is the total pages number and $d: ti \in d$ is the quantity of pages containing the term and log is based on 10.

$$tf = \frac{n_i}{n_{max}} \qquad (1)$$

$$IDF = \log\left(\frac{D}{d:t_i \in d}\right) \qquad (2)$$

### 3.2. Improved TF-IDF method

TF-IDF has the obvious deficiency if an important term appears in all the documents, so the denominator value will be small or zero according to IDF method [13]. For this reason, use the improved TF-IDF method to calculate the weight of each term. IDF defined avoid the zero in the result, therefore, it describes as illustrate as (3). To calculate the weight of each term, the improved weight formula is defined as (4), where, $w_i$ is representing the weight of of the term utilizing improved TF-IDF weighting strategy where for each term, and K=N-ni, N is the complete number of pages and ni is the quantity of pages that the term happens in.

$$IDF = log\left(\frac{D}{d:t_i \in d} + 0.01\right) \qquad (3)$$

$$w_i = \frac{TF_i * IDF_i * log(10 + \frac{n_i}{K}) * N}{\sqrt{\sum_{i=1}^{n}(TF_i * IDF_i)^2}} \qquad (4)$$

### 4. PROPOSED SYSTEM DESIGN

In system, users can search relevance social commerce web pages through the desired query as shown in Figure 1. An information retrieval system is utilized to retrieve information appropriate to a client's needs founded on inquiries presented to the framework by the client. The accompanying advances are proposed framework plan of this framework: To retrieve relevance web documents, a user poses a query to the system.
- This query is parsed by the system, and is then used to select relevance documents from the document collections
- Before parsing the query from the user, web document store often called to as the document collection in database, and these web documents extracted web contents from the database
- After that, the system performs two stages, namely, pre-processing stage and calculation of weights of terms
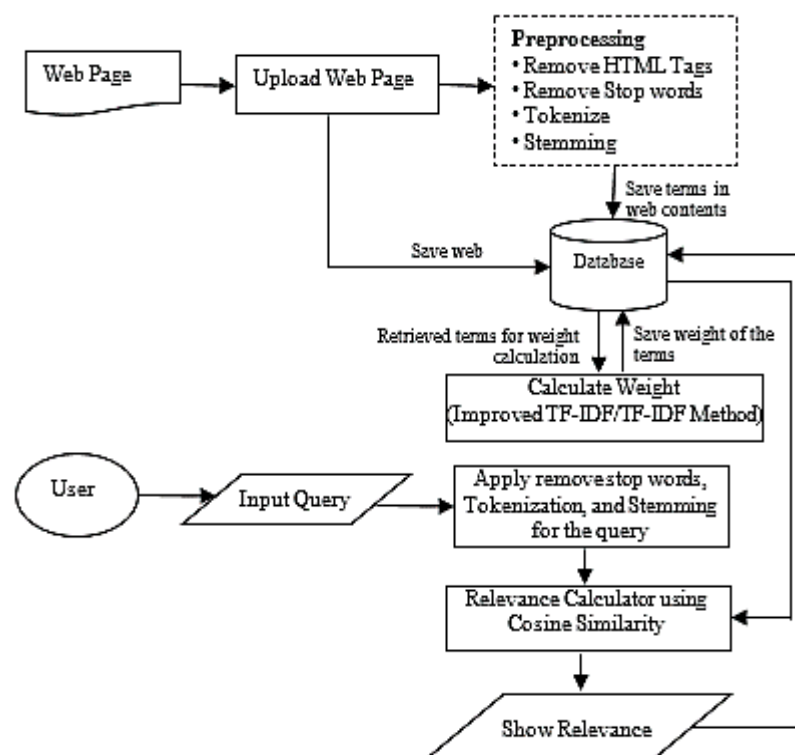


Figure 1. Proposed system design for comparison of improved TF-IDF method and TF-IDF method

In pre-processing stage: when hypertext markup language (HTML) documents income into the system, this document is tokenized. Tokenizing is to separate the blocks in this document, and remove the HTML tags. After tokenization, stop-words are removed. Stop-words include articles, pronouns, some of the verbs, nouns and adjectives. Poster stemming algorithm means a process for removing suffixes from words in English, to make text processing more efficient [14, 15]. This system is calculated weight of terms using (1) as TF, (2) as IDF and the product of TF and IDF in TF-IDF method. Also calculate weight of terms using (1), (3) and (4) of improved TF-IDF method [16]. For analyzing of comparison of improved TF-IDF method and IDF method based on weights of terms and then computes relevance calculation using cosine similarity method according to the results of relevance calculation and user's threshold value.

## 5.    TRUST EVALUATION

In this system, vector space model such as TF-IDF method and improved IF-IDF method determines the relevant web pages and irrelevant web pages using cosine similarity [17]. Let D={$OS_1$, $OS_2$..., $OS_M$} be a set of users reviews from online social commerce pages. OS means online social commence pages. Reviews from each page is uploaded to do preprocessing steps. The content data of five reviews are as following tables.

### 5.1.  Calculate with TF-IDF method

The first TF increases by IDF is TF-IDF, where IF and IDF are term frequency short term and inverse document frequency respectively [18]. Firstly, the numbers of IT words count in each review are as shown in Table 1. The system is calculated TF and IDF using (1) and (2) as shown Table 2. And then, the weights of terms are calculated using the product of TF and IDF results as shown in Table 3. And then calculate relevant pages based on user query using cosine similarity method. Cosine measure is to compute between document vector and *Qidf vector*.

$$Qidf = q_i * (idf)_I \tag{5}$$

Where $q_i$ means the number of words counts for each word in user comments and *idf* means inverse document query. The textual comparability between the page and the comments is the cosine similarity between *tf\*idf* vector of the query and the *tf\*idf* vector of each page [19, 20].

$$sim(Qidf, w_i) = \frac{\sum_{i=1}^{n} Qidf_i * w_i}{\sqrt{\sum_{i=1}^{n}(Qidf_i)^2} * \sqrt{\sum_{i=1}^{n}(w_i)^2}} \tag{6}$$

### Table 1. Number of word counts

| Term | OS1 | OS2 | OS3 | OS4 | OS5 |
|------|-----|-----|-----|-----|-----|
| Good | 0 | 1 | 0 | 0 | 4 |
| Great | 0 | 1 | 0 | 0 | 0 |
| Nice | 3 | 2 | 2 | 0 | 3 |
| Expensive | 0 | 1 | 1 | 1 | 0 |
| Satisfied | 0 | 1 | 1 | 2 | 1 |

### Table 2. Result of TF and IDF

| Term | OS1 | OS2 | OS3 | OS4 | OS5 | IDF |
|------|-----|-----|-----|-----|-----|-----|
| Good | 0 | 0.5 | 0 | 0 | 1 | 0.398 |
| Great | 0 | 0.5 | 0 | 0 | 0 | 0.6989 |
| Nice | 1 | 1 | 1 | 0 | 0.75 | 0.0969 |
| Expensive | 0 | 0.5 | 0.5 | 0.5 | 0 | 0.2219 |
| Satisfied | 0 | 0.5 | 0.5 | 1 | 0.25 | 0.0969 |

### Table 3. Weight of terms by TF-IDF

| Term | OS1 | OS2 | OS3 | OS4 | OS5 |
|------|-----|-----|-----|-----|-----|
| Good | 0 | 0.198 | 0 | 0 | 0.398 |
| Great | 0 | 0.349 | 0 | 0 | 0 |
| Nice | 0.097 | 0.097 | 0.097 | 0 | 0.073 |
| Expensive | 0 | 0.111 | 0.111 | 0.111 | 0 |
| Satisfied | 0 | 0.048 | 0.048 | 0.097 | 0.024 |

For calculation of cosine similarity method, compute *Qidf* vector based on user query using (5) as shown in Table 4. In this table, *Q* is the number of word counts in user query and *Qidf* is the product of word counts of each word and inverse document query values. And then, we calculate cosine similarity method using (6) as illustrate as Table 5. Therefore, all retrieved pages are rearranged to similarity values so that the relevant pages appear at top of the result set.

Table 4. Word counts and *Qidf* values

| Term | Q | Qidf |
|---|---|---|
| Good | 1 | 0.39794 |
| Great | 0 | 0 |
| Nice | 1 | 0.09691 |
| Expensive | 0 | 0 |
| Satisfied | 1 | 0.09691 |

Table 5. Result of cosine similarity method in TF-IDF

| OS1 | OS2 | OS3 | OS4 | OS5 |
|---|---|---|---|---|
| 0.230256 | 0.51413 | 0.215859 | 0.151493 | 0.983509 |

## 5.2. Improved TF-IDF method

TF values using (2) in improved TF-IDF method are the same values of TF-IDF method. The results of IDF values are calculated using (3) as shown in Table 6. The results of weight terms are calculated using (4) from values of TF and IDF in improved TF-IDF method. And then, the numbers of counts (Q) according to user query and calculate the results of *Qidf*. The cosine similarity method is calculated according to *Q* and *Qidf* based on user query as shown in Table 7. In this framework look at two strategies for unstructured inquiry positioning about their capacity to accurately disambiguate catch phrase question through systematic experiments. Document 5 is the top of the result set in relevant pages, but we analyzed that result value of cosine similarity method in improved TF-IDF method is less than the value in TF-IDF method.

Table 6. Results of improved IDF

| Term | OS1 | OS2 | OS3 | OS4 | OS5 | Improved IDF |
|---|---|---|---|---|---|---|
| Good | 0 | 0.5 | 0 | 0 | 1 | 0.399 |
| Great | 0 | 0.5 | 0 | 0 | 0 | 0.699 |
| Nice | 1 | 1 | 1 | 0 | 0.75 | 0.100 |
| Expensive | 0 | 0.5 | 0.5 | 0.5 | 0 | 0.224 |
| Satisfied | 0 | 0.5 | 0.5 | 1 | 0.25 | 0.100 |

Table 7. Results of cosine similarity method in improved TF-IDF

| OS1 | OS2 | OS3 | OS4 | OS5 |
|---|---|---|---|---|
| 0.23665 | 0.634956 | 0.258211 | 0.196793 | 0.97822 |

## 6. EXPERIMMENTAL RESULTS

In system experiments, to provide user flexible access and use relevant information, this system implements the results of two methods (improved TF-ID and TF-IDF) based on two domains such as bags social commerce and shoe social commerce. And then compare analysis of two methods using precision and recall as evaluation method [21, 22]. Recall is the ratio of the number of relevant topics retrieved to the total number of associate topics in database in (7). Precision is the proportion of the number of relevant topics retrieved to the irrelevant and associate with desire topics retrieved as (8).

$$Recall = \frac{No.of\ Relevant\ Documents\ Retireved}{Total\ No.of\ Documents\ Retrieved} \tag{7}$$

$$Precision = \frac{No.of\ Relevant\ Document\ Retrieved}{Total\ No.of\ Relevant\ Documents} \tag{8}$$

Experimental results include for both information retrieval with improved TF-IDF and TF-IDF methods to make performance analysis of social commerce domain among 50 usres reviews from 100 social commerce sites. Table 8 presents the experimental results with similarity threshold 0.1 for search using improved TF-IDF method [22]. When user searches information retrieval web pages using TF-IDF method produces better precision and recall as shown in Table 9.

Table 8. Experimental results of improved TF-IDF method (bags social commerce)

| Term | Relevance Pages | Particular Pages in DB | Total Retrieved Pages | Precision | Recall |
|---|---|---|---|---|---|
| Good | 5 | 32 | 5 | 0.7650 | 0.1562 |
| Great | 4 | 53 | 4 | 0.8059 | 0.7755 |
| Nice | 4 | 38 | 4 | 0.6890 | 0.6053 |
| Expensive | 11 | 40 | 11 | 0.8930 | 0.8275 |
| Satisfied | 6 | 28 | 6 | 0.7532 | 0.6143 |

Table 9. Experimental results of TF-IDF method (bags social commerce)

| Term | Relevance Pages | Particular Pages in DB | Total Retrieved Pages | Precision | Recall |
|---|---|---|---|---|---|
| Good | 31 | 32 | 36 | 0.8611 | 0.9688 |
| Great | 44 | 53 | 51 | 0.8627 | 0.8302 |
| Nice | 33 | 38 | 37 | 0.8919 | 0.8684 |
| Expensive | 31 | 40 | 42 | 0.7381 | 0.7751 |
| Satisfied | 25 | 28 | 30 | 0.8333 | 0.8929 |

Thus, in bags social commerce domain, we can see TF-IDF method is the best recall than the improved TF-IDF according to experimental result. In our experimental, we use 0.1 as threshold value for both TF-IDF method and improved TF-IDF method. When threshold value is higher (<0.1), recall and precision result is not good in improved TF-IDF method in Tables 10 and 11. In testing can see TF-IDF method is best recall than the improved TF-IDF method. Hence in implementation, use 0.01 threshold value for both TF-IDF method and improved TF-IDF method. When threshold value is higher, recall and precision result is not good in improved TF-IDF method. Based on the above results considerations, this paper's experiments in the N value of 120, as a result of 5 social commerce feature words can described the subject information of a selected web pages.

Table 10. Experimental results of improved TF-IDF method (shoes social commerce) (similarity threshold 0.01)

| Term | Relevance Pages | Particular Pages in DB | Total Retrieved Pages | Precision | Recall |
|---|---|---|---|---|---|
| Good | 38 | 41 | 39 | 0.9744 | 0.9268 |
| Great | 45 | 50 | 48 | 0.9375 | 0.9 |
| Nice | 40 | 43 | 50 | 0.8677 | 0.9302 |
| Color | 40 | 45 | 43 | 0.9302 | 0.8889 |
| Size | 37 | 44 | 47 | 0.7872 | 0.8409 |

Table 11. Experimental results of TF-IDF method (shoes social commerce) (similarity threshold 0.01)

| Term | Relevance Pages | Particular Pages in DB | Total Retrieved Pages | Precision | Recall |
|---|---|---|---|---|---|
| Good | 39 | 41 | 47 | 0.8298 | 0.9512 |
| Great | 48 | 50 | 57 | 0.8421 | 0.9605 |
| Nice | 40 | 43 | 68 | 0.5882 | 0.9302 |
| Color | 40 | 45 | 51 | 0.7843 | 0.8889 |
| Size | 25 | 47 | 30 | 0.8333 | 0.5319 |

The recall ratio and precision of TF-IDF algorithm and improved TF-IDF algorithm are shown in Figure 2 and Figure 3. After the calculation the result can be seen from the Figure 2, the recall rate of the improved TF-IDF algorithm. Also, can be seen from the Figure 3, the precision of the improved TF-IDF-compared with the classic TF-IDF algorithm. That can be seen from the test results, in similar conditions, utilizing the improved equation of overall performance is simple implementation, easy to understand calculation and strong explanatory that better than utilizing TF-IDF method [23-25].
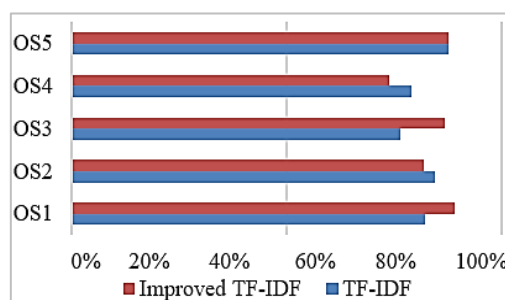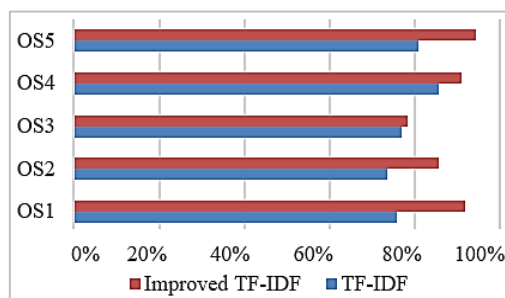


Figure 2. Recall ratio

Figure 3. Precision

## 7. CONCLUSIONS

With the quick development of data on the World Wide Web, finding and retrieving valuable data turns into a significant issue. A web client may use the positioned site page list for exploring the web and finding pertinent pages. The point of this framework computes values from loads of the terms to give expectation of significant pages (social commerce web sites). This system examines relevance web pages by the comparison of TF-IDF method and Improved TF-IDF method. We conclude that nearly all relevant web documents based on user query extracted by using TF-IDF method. Likewise, TF-IDF accomplishment in foreseeing trust for user dependent on remarks in shares by post and furthermore vector space model is simpler to find the most confided in way as least time more conceivable than other and at the time of trust computation. Later on, user's audits are totally significant and find out progressively proper procedures dependent on this framework for trust in online communities.

## REFERENCES

[1] Lu S., Li X., Bai S., Wang S., "An improved approach to weighting terms in text," *Journal of Chinese Information Processing,* vol. 14, no. 6, pp. 8-13, 2000.
[2] Mohammed Zuhair Al-Taie, Seifedine Kadry, Joel Pinho Lucas, "Online Data Preprocessing: A Case Study Approach," *International Journal of Electrical and Computer Engineering (IJECE),* vol. 9, no. 4, 2620-2626, 2019.
[3] Hla Sann Sint, Khine Khine Oo., "Consumer Trust Recommendation in Online Social Commerce," *2019 IEEE 8th Global Conference on Consumer Electronics (GCCE),* pp. 445-447, 2019.
[4] Hassan Najadat, Amer Al-Badarneh, Sawsan Alodibat, "A Review of Website Evaluation Using Web Diagnostic Tools and Data Environment Analysis," *Bulletin of Electrical Engineering and Informatics*, vol. 10, no 1, pp. 258-265, 2021.
[5] Darshana Karna, Ilsang Ko, "We-Intention, Moral Trust and Self-Motivation on Accelerating Knowledge Sharing in Social Collaboration," *48th Hawaii International Conference on System Sciences*, 2015.
[6] Xiaoming Zheng, Yan Wang, Mehmet A. Orgun, "BiNet: Trust Sub - network Extraction using Binary Ant Colony Algorithm in Contexual Social networks," *2015 IEEE International Conference on Web Services,* 2015.
[7] Ye Hanmin, Zhang Qiuling, Huang Peiliang, "Collaborative filtering algorithm Research based on Matrix Factorization and Multi-path Trust Degree Fusion," *3rd International Conference on Information Science and Control Engineering (ICISCE)*, 2016.
[8] Zhang C., Wang H., Liu Y., and Xu H., "Document Clustering Description Extraction and Its Application," *Proceedings of the 22nd International Conference on Computer Processing of Oriental Languages,* 2015, pp. 370-377.
[9] Ma Zhanguo, Jing Feng, Liang Chen, Xiangyi Hu, and Yanqin Shi. "An Improved Approach to Terms Weighting in Text Classification," *2011 International Conference on Computer and Management (CAMAN),* 2011.
[10] Ahed M. F. Al Sbou, "A survey of Arabic text classification models," *International Journal of Informatics and Communication Technology (IJ-ICT),* vol. 8, no. 6, pp. 4352-4355, 2019.
[11] Ali Hameed Yassir, Ali A. Mohammed, Adel Abdul-Jabbar Alkhazraji, Mustafa Emad Hameed, Mohammed Saad Talib, Mohanad Faeq Ali, "Sentimental classification analysis of polarity multi-view textual data using data mining techniques," *International Journal of Electrical and Computer Engineering (IJECE),* vol. 10, no. 5, pp. 5526-5534, 2020.
[12] Peter Hakansson, Hope Witmer, "Social Media and Trust – A Systematic Literature Review," *Journal of business and economics*, vol. 6, no. 3, pp. 517-524, 2015.
[13] Xiaohui Cui, Thomas E. Potok, "Chapter 10 A Distributed Agent Implementation of Multiple Species Flocking Model for Document Partitioning Clustering," *Springer Science and Business Media LLC*, 2006.
[14] Bo Yang, Yu Lei, Jiming Liu, Fellow, "Social Collaborative Filtering by Trust," *Transactions on Pattern Analysis and Machine Intelligence,* vol. 39, no. 8, pp. 1633-1647, 2017.
[15] Edi Sutoyo, Ahmad Almaarif, "Twitter sentiment analysis of the relocation of Indonesia's capital city," *Bulletin of Electrical Engineering and Informatics,* vol. 9, no. 4, pp. 1620-1630, 2020.
[16] Miguel E. Ruiz. "Evaluating topic-driven web crawlers," *Proceedings of the 24th annual international ACM SIGIR conference on Research and development in information retrieval,* 2001.
[17] Qi Fu, Jun Tan. "Chapter 41 Research on Detection and Trend Forecasting Technologies of Micro-blog Hot Topic," *Springer Science and Business Media LLC*, 2017.

[18] F. Xhafa, S. Patnaik, A. Y. Zomaya, "Advances in Intelligent Systems and Interactive Applications," *Springer Science and Business Media LLC*, 2020.

[19] Heru Agus Santoso, *et al.,* "Hoax classification and sentiment analysis of Indonesian news using Naive Bayes optimization," *TELKOMNIKA Telecommunication Computing Electronics and Control,* vol. 18, no. 2, pp. 799-806, 2020.

[20] Huajia Wang, Ruo Hu, Hong Xu, "Research on an Improved Algorithm of Professional Information Retrieval System," *Proceedings of the 2019 3rd International Conference on Digital Signal Processing - ICDSP 2019*, 2019.

[21] Bing Zhu, Yang Yu, Chuanzhen Li, Hui Wang, "Research and implementation of hot topic detection system based on web," *2017 3rd IEEE International Conference on Computer and Communications (ICCC)*, 2017.

[22] Guoqin Chang, Hua Huo, "A Method of Fine-Grained Short Text Sentiment Analysis Based on Machine Learning," *Neural Network World*, vol. 28, no. 4, pp. 345-360, 2018.

[23] Ayad, H. and Kamel, M. S., "Topic Discovery from Text Using Aggregation of Different Clustering Methods," *15th Conference of the Canadian Society for Computational Studiesof Intelligence on Advances in Artificial Intelligence*, 2002, pp. 161-175.

[24] Mimi Aminah binti Wan Nordin, *et al.,* "The disruptometer: an artificial intelligence algorithm for market insights," *Bulletin of Electrical Engineering and Informatics*, vol. 8, no. 2, pp. 727-734, 2019.

[25] Amany A. Naem, Neveen I. Ghali, "Optimizing community detection in social networks using antlion and K-median," *Bulletin of Electrical Engineering and Informatics*, vol. 8, no. 4, pp. 1433-1440, 2019.

## BIOGRAPHIES OF AUTHORS

**Hla Sann Sint** is a PhD candidate at University of Computer Studies, Yangon. She is very interested in information retrieval, knowledge based trust, opinion mining and social networking. Her current research is related to knowledge based trust, opinion mining and social networking. She is also a research student at University of Computer Studies, Yangon, now.



**Professor Dr. Khine Khine Oo** is a head of Software Engineering Lab at University of Computer Studies, Yangon (UCSY). She is interested in Software Engineering, information retrieval and opinion mining. Her research areas include software engineering, opinion mining, database, social networking and information retrieval. She currently serves at the Faculty of Information Science, University of Computer Studies, Myanmar.