

Comparative Analysis of TF-IDF and TF-IDF Groups using BERT and GloVe

Alex Conro Manuel¹[12S18022] and Angelina Naomi Christina Sinaga¹[12S18035]

Information System, Del Institute of Technology, Indonesia
{iss18022, iss18035}@students.del.ac.id

Abstract. Rigidity in finding information and a lot of noise are problems that often occur in knowledge management which has the concept of information retrieval in identifying, explaining and distributing information for use obtained from documents or collections. Currently the development of information retrieval is very interesting to discuss and research, because the application of information retrieval can help overcome some of the problems above. The use of TF-IDF has been widely used because it is simple in the process of calculating keywords or queries and is easy to use in measuring content uniqueness, as well as low-cost computational processes. However, so far the implementation of the TF-IDF group on IR has not been implemented. In this project, we compare the performance of TF-IDF with TF-IDF Group by using word embedding method: BERT and GloVe, where the document grouping process is done using Minibatch K-Means (Cosine Sim). The first stage is text preprocessing which consists of case folding, stopword, tokenization and stemming stages, the second stage is weighting using TF-IDF and TF-IDF Group, the third stage is applying the method used and the last stage is evaluating.

Keywords: information retrieval · BERT · GloVe.

1 Introduction

Information dissemination is currently growing rapidly. Every day people search for information by typing keywords in search engines and want fast and accurate information [1]. Information consists of various categories and is scattered randomly and unclearly [19], which is sometimes difficult for those who need it, moreover, its contents do not necessarily contain important things needed by readers. This information search activity is known as information searching. Information search aims to find the most relevant documents based on keywords in user-generated queries [13]. Basically, the development of this information retrieval system is actually inseparable from the techniques or methods used. There are two jobs in this system, namely preprocessing the dataset and applying certain methods to calculate the relevance (similarity) between documents in the preprocessed database [9]. As a result, the system will return a list of documents ordered according to the similarity value to the previously entered query. The solution to this problem is to summarize the text [20].

In this paper, we analyze the comparison of the TF-IDF Group with the TF-IDF using BERT and GloVe. Considering the TF-IDF group and TF-IDF, this comparative analysis aims to compare the output to the performance of text categorization. The TF-IDF weighting consists of 2 factors, namely term frequency (TF) and inverse document frequency (IDF). Term frequency (TF) is a condition in which each term is assumed to have a proportion of importance according to the number of occurrences in the document [18] and inverse document frequency (IDF) is a term weighting method that focuses on paying attention to the occurrence of terms in the entire text collection [15]. There are several parameters that are used as benchmarks to compare the performance of the text categorization, namely precision, recall and f-measure [3]. Minibatch K-means is a version of the standard K-means algorithm in machine learning that uses small, random, and fixed-sized data sets to be stored in memory, and then with each iteration, a random sample of data is collected and used to update the cluster [10]. Global Vectors for Word Representation (GloVe) is a word representation to generate word embedding to be used to handle word similarity, word analogy, and named entity recognition [6]. Transformers' Bidirectional Encoder Representations (BERT) is a neural network-based technique for pre-training natural language to help understand the context of words in search queries [12]. Based on the above, we first need to combine two concepts for calculation, namely the frequency of occurrence of a word in a particular document and the inverse of the frequency of documents containing that word against the BERT method and the Glove method. Second, calculate the TF-IDF group against the BERT method and the Glove method. Finally, analyze the results from the first and second stages, then make a comparison.

Based on the research of Kamyab et al. [8] proposed a new attention-based model that utilizes CNN with LSTM (named ACL-SA), applies a preprocessor to improve data quality and uses term frequency-inverse document frequency (TF-IDF) feature weighting and Glove's pre-trained word embedding approach to extract meaningful information from textual data, use CNN max-pooling to extract contextual features and reduce feature dimensions, also use integrated two-way LSTM to capture long-term dependencies. In the research of Weilong Chen et al. [4], it focuses on the effect of different contexts to determine the similarity of 2 different words. Their research is based on BERT built with TF-IDF and applies the data collection method (CoSimLex), which covers four languages namely English, Croatian, Slovenian and Finnish. In the model they built word embedding can train the model to predict the similarity of words to understand the meaning of words from different perspectives. Research Jin et al. [7] created a multi-label classification framework for aspect-based sentiment analysis problems in restaurant customer reviews where their processes include text preprocessing, feature extraction using modified BERT and TF-IDF, and fine tuning. The TF-IDF method is used to determine how important the word is in the multi-label classification by calculating the weights. [17].

Although the above studies have proven the superiority of each method in calculating the weights, none of them compared the weighting of the TF-IDF

with the TF-IDF Group using BERT and GloVe. Based on this, we propose a Comparative Analysis of TF-IDF with TF-IDF Group using BERT and GloVe to find out what is the special differentiator in calculating the weights. Specifically, we will do text preprocessing which consists of case folding, stopwords, tokenization and stemming stages. Then we apply the calculation of the frequency of occurrence of a word in a particular document and inverse the frequency of documents containing the searched word using both methods. Then determine the ranking using cosine similarity and followed by the final analysis of the results obtained in the BERT and GloVe methods. The main contributions are as follows:

1. We propose the BERT and GloVe models to make comparisons on the TF-IDF and TF-IDF Group.
2. We Group documents using Minibatch K-Means (Cosine Sim). Before grouping the document, do docs expansion (following the concept of query expansion)
3. We conducted experiments on two sets of datasets, namely Spam and BBC News to compare the results of the two applied methods.

The rest of the paper is organized as follows: Section 2 discusses the related works. Then, we present the architecture of the framework in Section 3. Next, Section 4 provides the experiment setup and implementation and experimental results and analysis, and finally, Section 5 concludes this work.

2 Related Works

We first review the important things related to what we are working on, namely the expansion document, Bert, GloVe, TF-IDF and TF-IDF Group.

Document Expansion: Document expansion is method using neural networks beside it document expansion more effective than query expansion on these two datasets most likely because there are more signals to exploit as documents are much longer [14]. Based on Rocchio relevance feedback where D is there a document to expand, $\left(D \frac{N}{j}\right) \frac{N}{j} = 1$, is the set of R neighborhood documents with D as the query, and D' is the extended document.

$$D' = \alpha D + \frac{1 - \alpha}{R} \sum_{j=1}^R D \frac{N}{j} \quad (1)$$

Equation 1 reweights the original term in the current document by a factor, and reweights the additional term from the neighborhood of D by a factor $(1 - \alpha)/R$. Note that every document $D \frac{N}{j}$ from neighborhood D in Equation 1 contributes equally to the expansion. However, a method for re-weighting terms in documents from the neighborhood by a factor proportional to their similarity to the neighborhood is proposed in, as in Equation 2. the proportionality factor for documents $D \frac{N}{j}$ the neighborhood is the similarity ratio. with D , with total similarity for all documents in the neighborhood. And hypothesizes whether documents in

that environment are fetched higher rankings and are more similar to current documents and thus more reliable for expansion.

$$D' = \alpha D + (1 - \alpha) \sum_{j=1}^R \frac{\text{sim}(D, D_{\frac{N}{j}})}{\sum_{k=1}^R \text{sim}(D, D_{\frac{N}{k}})} D_{\frac{N}{j}} \quad (2)$$

The intuitive idea of using proportional equality weights for various documents is presented theoretically by the RLM model, by calculating the relevance model $P(w|R)$, in the document expansion RLM estimates a new document model created with the symbol D' and returns the current D document and the neighborhood $\left(D_{\frac{N}{j}}\right) \frac{N}{j} = 1$, assuming that, *assuming that*, and the probability $P(w|D')$, is estimated by $P(w, d_1, \dots, d_n)$, is given by

$$P(w; d_1, \dots, d_n) = \sum_{j=1}^R P\left(D_{\frac{N}{j}}\right) P(w; d_1, \dots, d_n | D_{\frac{N}{j}}) = \frac{1}{R} \sum_{j=1}^R P(w | D_{\frac{N}{j}}) \prod_{i=1}^n P(d_i | D_{\frac{N}{j}}) \approx P(w | D') \quad (3)$$

The language model for the document expansion D' is used during the retrieve, recompute the estimation of the relevance model $P(w|D)$ or the expansion such as equation 3 and the unigram document model for D , which is shown in equation 4 as follows:

$$P(w | D') = \alpha P(w | D) + (1 - \alpha) P(w | D') = \alpha P(w | D) + \frac{1 - \alpha}{R} \sum_{j=1}^R P(w | D_{\frac{N}{j}}) \prod_{i=1}^n P(d_i | D_{\frac{N}{j}}) \quad (4)$$

Note that Equation 4 is similar to Equation 2 in that the quantity $\prod_{i=1}^n P(d_i | D_{\frac{N}{j}})$, Which acts as the proportional similarity of the neighborhood document $D_{\frac{N}{j}}$ to D . As how DE works, it consists of four steps [2]:

1. Preprocessing original terms and additional terms
2. Ranking
3. Define a cost function.
4. Document reformulation.

BERT: BERT uses a multi-layer network to capture and Transformers to encode input tokens, and make them a representative model, the token output vector becomes a vector that is used to represent the token sentence [5]. **Cosine Similarity:** The Cosine similarity score represents the scenario that the user input exactly matches the document; the semantic score represents the scenario that the user wants to search for some relevant document [11]. **GloVe:** GloVe is a method that combines co-occurrence and semantic relationships and is a global matrix factorization method, which represents the occurrence of a word in a document. Where GloVe studies the relationship of a word by calculating the frequency of occurrence of the word along with other words in a given corpus. This frequency of occurrence has the potential to encode multiple forms of pronouns and help performance in word analogies. GloVe stages are:

1. Collect word co-occurrence statistics in the form of a word co-occurrence matrix.
2. Define soft constraints on word pairs. w_i is the main vector, w_j is the context vector, b_i , b_j is the scalar bias for the main and context words.

$$W \frac{T}{i} \frac{T}{w_j} + b_i + b_j = \log(x_{ij}) \quad (5)$$

3. Define a cost function.

$$J = \sum_{i=1}^v \sum_{j=1}^v f(x_{ij}) W \frac{T}{i} \frac{T}{w_j} + b_i + b_j = \log(x_{ij})^2 \quad (6)$$

f is a weighting function to help prevent the learning of common word pairs. The function is defined as follows:

$$f(x_{ij}) = \left\{ \left(\frac{x_{ij}}{x_{max1}} \right) \text{ if } x_{ij} < XMAX \right\} \text{ otherwise} \quad (7)$$

TF-IDF: The TF-IDF models the TF-IDF Weighted score. using Weighted TF-IDF to rank documents, in both Glove and BERT Glove methods we calculate the similarity between these user and candidate queries, then re-rank the documents and get similar TF-IDF scores. TF-IDF gives the same score regardless of different semantic information, we aim to compare the results of both methods and datasets. TF-IDF is used to calculate the relevance of a word in a particular document by multiplying two matrices between the frequency of words in a document with the frequency of word documents in a set of documents [16], calculated as follows:

$$tfidf(t, d, D) = tf(t, d), idf(t, D) \quad (8)$$

$$tf(t, d) = \log(1 + freq(t, d)) \quad (9)$$

$$idf(t, D) = \log\left(\frac{N}{count(d \in D : t \in d)}\right) \quad (10)$$

3 Framework

In this section, we search for TF-IDF values and TF-IDF Group values using the BERT and GloVe methods, where the document grouping process is carried out using Minibatch K-Means (Cosine Sim), and an overview is shown in Figure 1 and Figure 2.

The experiment that will be carried out starts with data preprocessing which consists of case folding, tokenization, removing stop words and stemming. After preprocessing the data, a clean dataset will be obtained. This clean dataset will then be clustered using mini batches, so that clusters are obtained based on sentences. After the clustering process continues with the split dataset. From the split dataset, word embedding will be performed

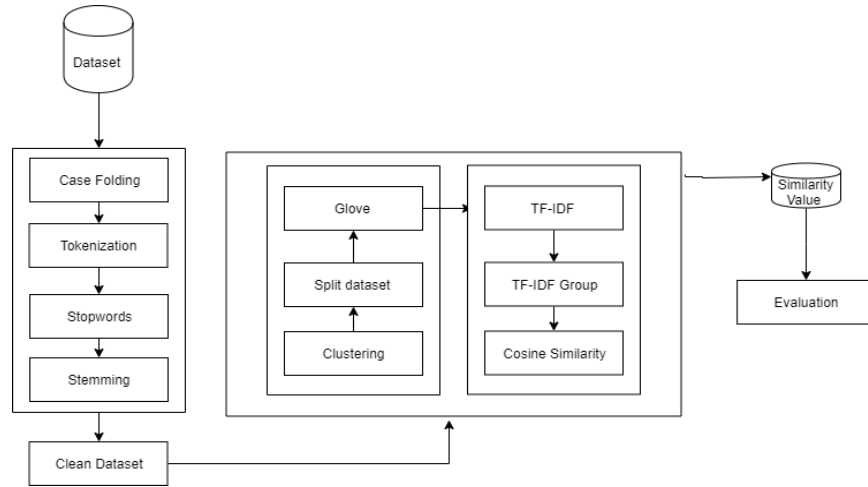


Fig. 1: Framework Glove

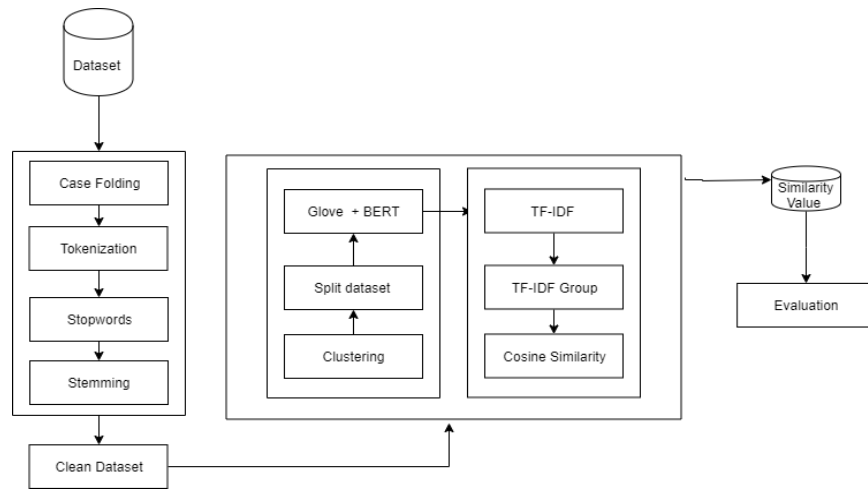


Fig. 2: Framework GloveBERT

with Glove for figure 1, or after word embedding is inputted into the BERT model as in figure 2. Next, a weighting scheme will be carried out using TF-ID and checking for similarity with cosine similarity. To get this similarity value, a search will be carried out on one of the documents. After that, the searched document will be obtained and an evaluation will be carried out, which method is good or not.

4 Experimental Evaluation

In this section, we discuss the evaluation of BERT and Glove by conducting comprehensive experiments on two datasets namely the spam dataset and the BBC News Dataset and show which model performs better on which dataset..

4.1 Experimental Setup

We conducted experiments on two publicly available datasets, namely the Spam Dataset and the BBC News Dataset, which are the datasets commonly used for information retrieval. The Spam Dataset is a collection of SMS tagged messages that have been collected for SMS Spam research. It contains a set of SMS messages in English from 5,574 messages, marked as ham (legitimate) or spam. While the BBC News dataset is a collection of RSS feeds from BBC News which consists of several columns such as title, date and description.

Table 1: Dataset

Dataset	Items	Size(KB)
Spam	5574	492
BBC News	1816	585

To find out how the performance of the model, it is proposed to calculate with accuracy. Accuracy is defined as the level of closeness between the predicted value and the actual value.

4.2 Evaluation Methodology

The results from the system certainly need to be measured. Evaluation measurements that are very commonly used as a reference in determining the effectiveness of information. Precision is match between the request information and the results from request, which depends on relevance, where relevance is how appropriate the document is for the purposes of seeking information, and how relevant the document is to the seeker. To calculate the accuracy use the following formula:

$$Accuracy(Acc) = \frac{TP + TN}{TP + TN + FP + FN} \quad (11)$$

4.3 Evaluation on Effectiveness

The effectiveness comparison that we compare to Glove method and the Glove BERT by performing the same two dataset ,to the same target set. Table describes the comparison results of Glove and Glove BERT,in terms of Cosine similarity, and accuracy.

Table 2: Effectiveness comparison with for *Spam and BBC* Dataset in terms of cosine similarity and accuracy (the lower the value, the better).

Method	Dataset	Cosine Similarity	Accuracy
Glove	Spam	0.710	0.977
Glove BERT	Spam	0.677	0.954
Glove	BBC News	0.683	0.058
Glove BERT	BBC News	0.640	0.0010

Configuration Parameter. In this study, we did not configure parameters, we used parameters and their values are default values. Based on table 3 above, the Glove method has cosine similarity value and higher accuracy than the Glove-BERT method in the second dataset. Apart from the accuracy value between the Spam and BBC News datasets which are quite far apart, it has nothing to do with methods and configurations, this is only a dataset problem.

5 Conclusion

In this paper, we propose a comparison of the TF-IDF with the TF-IDF Group using Glove method, from the results of the experiments we carried out on two datasets, the cosine similarity and accuracy were high. Merging Glove with BERT actually gives a bad effect, the addition of BERT actually make value of cosine similarity and accuracy is low, besides that it is necessary to check what dataset is used, so that this does not happen like in this paper where the accuracy value in the second dataset is much different from the dataset one.

References

1. https://balitbangsdm.kominfo.go.id/upt/bandung/?mod=publikasi&action=dl&cid=10&pub_id=5
2. Azad, H.K., Deepak, A.: Query expansion techniques for information retrieval: a survey. *Information Processing & Management* **56**(5), 1698–1735 (2019)
3. Béjar Alonso, J.: K-means vs mini batch k-means: A comparison (2013)
4. Chen, W., Yuan, X., Zhang, S., Wu, J., Zhang, Y., Wang, Y.: Ferryman at semeval-2020 task 3: bert with tfidf-weighting for predicting the effect of context in word similarity. In: *Proceedings of the fourteenth workshop on semantic evaluation*. pp. 281–285 (2020)
5. Fimoza, D., et al.: Analisis sentimen terhadap film indonesia dengan pendekatan bert (2021)
6. Fitri, M.: perancangan sistem temu balik informasi dengan metode pembobotan kombinasi tf-idf untuk pencarian dokumen berbahasa indonesia. *JUSTIN (Jurnal Sistem dan Teknologi Informasi)* **1**(1), 80–85 (2013)
7. Jin, Z., Lai, X., Cao, J.: Multi-label sentiment analysis base on bert with modified tf-idf. In: *2020 IEEE International Symposium on Product Compliance Engineering-Asia (ISPCE-CN)*. pp. 1–6. IEEE (2020)
8. Kamyab, M., Liu, G., Adjeisah, M.: Attention-based cnn and bi-lstm model based on tf-idf and glove word embedding for sentiment analysis. *Applied Sciences* **11**(23), 11255 (2021)
9. Kesuma, H.W.A.: Penerapan metode tf-idf dan cosine similarity dalam aplikasi kitab undang-undang hukum dagang (2016)
10. Khatri, A., et al.: Sarcasm detection in tweets with bert and glove embeddings. *arXiv preprint arXiv:2006.11512* (2020)
11. Marwah, D., Beel, J.: Term-recency for tf-idf, bm25 and use term weighting. In: *Proceedings of the 8th International Workshop on Mining Scientific Publications*. pp. 36–41 (2020)
12. Melita, R.: Penerapan Metode Term Frequency Inverse Document Frequency (Tf-Idf) Dan Cosine Similarity Pada Sistem Temu Kembali Informasi Untuk Mengetahui Syarah Hadits Berbasis Web (Studi Kasus: Hadits Shahih Bukhari-Muslim). B.S. thesis, Fakultas Sains dan Teknologi UIN Syarif Hidayatullah Jakarta (2018)
13. Merdekawan, R., et al.: Sistem Penelusuran Katalog Perpustakaan Menggunakan Metode Rocchio Relevance Feedback. Ph.D. thesis, Universitas Kanjuruhan Malang (2014)
14. Nogueira, R., Yang, W., Lin, J., Cho, K.: Document expansion by query prediction. *arXiv preprint arXiv:1904.08375* (2019)
15. Putra, J.W.G.: Pengenalan konsep pembelajaran mesin dan deep learning. Tokyo. Jepang (2019)
16. Qaiser, S., Ali, R.: Text mining: use of tf-idf to examine the relevance of words to documents. *International Journal of Computer Applications* **181**(1), 25–29 (2018)
17. Rachman, F.P., Santoso, H., Djajadi, A.: Machine learning mini batch k-means and business intelligence utilization for credit card customer segmentation. *International Journal of Advanced Computer Science and Applications* **12**(10) (2021)
18. Saadah, M.N., Atmagi, R.W., Rahayu, D.S., Arifin, A.Z.: Sistem temu kembali dokumen teks dengan pembobotan tf-idf dan lcs. *Jurnal Ilmiah Teknologi Informasi* **11**(1), 19–22 (2013)
19. Widyasanti, N.K., Putra, I.K.G.D., Rusjyanthi, N.K.D.: Seleksi fitur bobot kata dengan metode tfidf untuk ringkasan bahasa indonesia. *Jurnal Ilmiah Merpati (Menara Penelitian Akademika Teknologi Informasi)* pp. 119–126 (2018)

20. Winarti, T., Setiawan, W., Iswoyo, Pudjiastuti, E.: Deteksi kemiripan dokumen bahasa indonesia dengan menggunakan model ruang vektor, <http://repository.usm.ac.id/files/research/G001/20191127114354-Deteksi-Kemiripan-Dokumen-Bahasa-Indonesia-Dengan-Menggunakan-Model-Ruang-Vektor.pdf>

6 Contribution

This section describes the division of tasks in this project.

Table 3: Contribution

Nama	Document	Dataset with Approach	Code	log(21 Mei)	log(22 Mei)
Alex Conro Manuel	Related Works, Experimental Setup, Evaluation Methodology, Evaluation on Effectiveness, conclusion	Spam ,Glove	EDA, Data Pre-processing, Feature Extraction (TF-IDF), TF-IDF, Word Embeddings: GloVe, K-means(mini batch), TF-IDF VS TF-IDF Group	Add evaluation in code, Update Evaluation Methodology	add cosine simil
Angelina Naomi Christina Sinaga	Abstract, Introduction, Framework	BBC News, Glove	EDA, Data Pre-processing, Feature Extraction (TF-IDF), TF-IDF, Word Embeddings: GloVe		

Link gitbub : <https://github.com/angelinasinaga/Comparative-Analysis-of-TF-IDF-and-TF-IDF-Groups-using-BERT-and-GloVe>

Table 4: Contribution

Nama	Document	Dataset with Approach	Code	log(21 Mei)	log(22 Mei)
Alex Conro Manuel	Related Works, Experimental Setup, Evaluation Methodology, Evaluation on Effectiveness, conclusion	Spam ,BERT-Glove	EDA, Data Pre-processing, Feature Extraction (TF-IDF), TF-IDF, Word Embeddings: GloVe, K-means(mini batch), add visualization, add TF-IDF VS TF-IDF Group, BERT	Add evaluation in code	add cosine simil
Alex Conro Manuel	Related Works, Experimental Setup, Evaluation Methodology, Evaluation on Effectiveness, conclusion	BBC ,BERT-Glove	EDA, Data Pre-processing, Feature Extraction (TF-IDF), TF-IDF, Word Embeddings: GloVe, K-means(mini batch), TF-IDF VS TF-IDF Group, BERT		add evaluation, cosine similarity